

# Networks

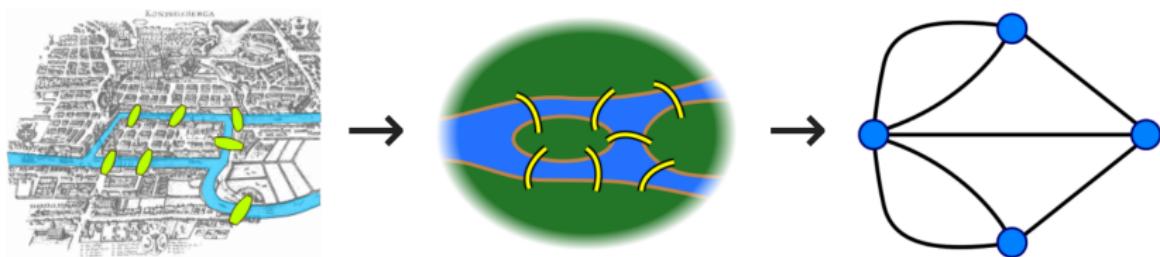
Arthur Charpentier

UQAM

Actuarial Summer School 2019

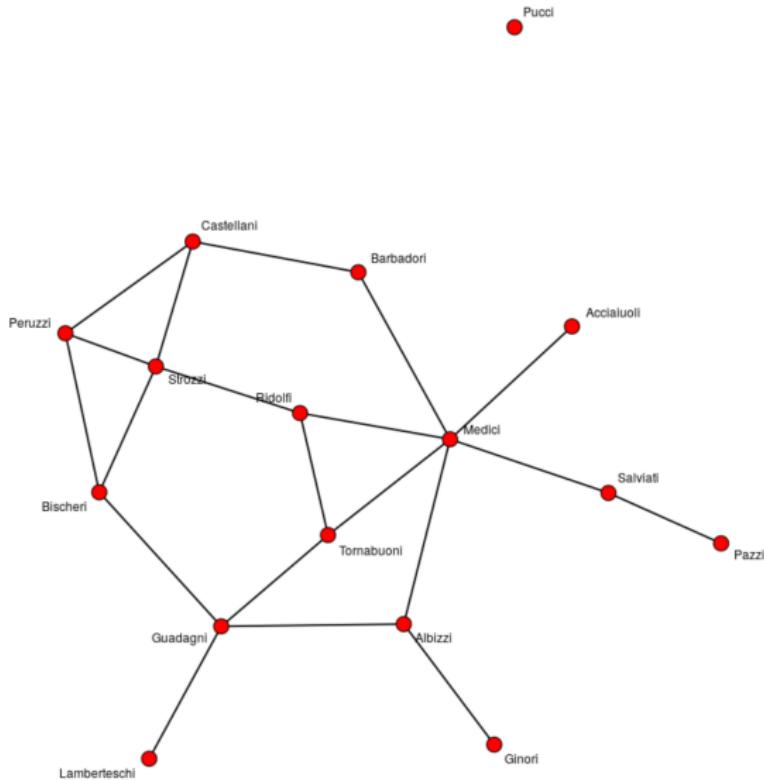
# Königsberg's Seven Bridges

The path should cross over each of the seven bridges exactly once.



see also use of GPS to find the shortest path home

# Florentine Marriages (1430's)



## References

Briatte's [Awesome Network Analysis](#), Adamic (2010, [lecture notes](#)), Easley & Kleinberg (2010, [Networks, Crowds, and Markets: Reasoning About a Highly Connected World](#)), Harary (1969, [Graph Theory](#)), Jackson (2008, [Social and Economic Networks](#)), Kolaczyk (2009, [Statistical Analysis of Network Data: Methods and Models](#)), Kolaczyk & Csárdi (2010, [Statistical Analysis of Network Data with R](#)) Newman (2004, [The structure and function of complex networks](#)), Newman (2010, [Networks : An Introduction](#)), Wasserman & Faust (1994) [Social Network Analysis: Methods and Applications](#)), West (1996). [Introduction to Graph Theory](#)

See also R packages, [sna](#), [network](#), [igraph](#) to create and manipulate networks, and [ggnet2](#) or [networkD3](#).

# Network Representation

Let  $V = \{1, \dots, n_V\}$  denote either **nodes**, or **vertices** ( $n_V$  is the **order**). Let  $E \in \{0, 1\}^{n_V \times n_V}$  represents the **relationships**, through an adjacency matrix  $A$ ,  $A_{i,j} = 1$  indicates a **link** - or **edge** - between  $i$  and  $j$ , or a collection of links  $\{e_1, \dots, e_{n_E}\}$ . Let  $n_E = |E|$  denote the number of edges, called **size**.

The **degree**  $d(\cdot)$  of a vertice  $v$  is its number of incident edges.  
A **network** is a pair  $G = (V, E)$

$G = \text{Internet}$ ,  $V = \text{computers}$ ,  $E = \text{IP network adjacency}$

$G = \text{World Wide Web}$ ,  $V = \text{web pages}$ ,  $E = \text{hyperlink}$

$G = \text{Articles}$ ,  $V = \text{authors}$ ,  $E = \text{citations}$

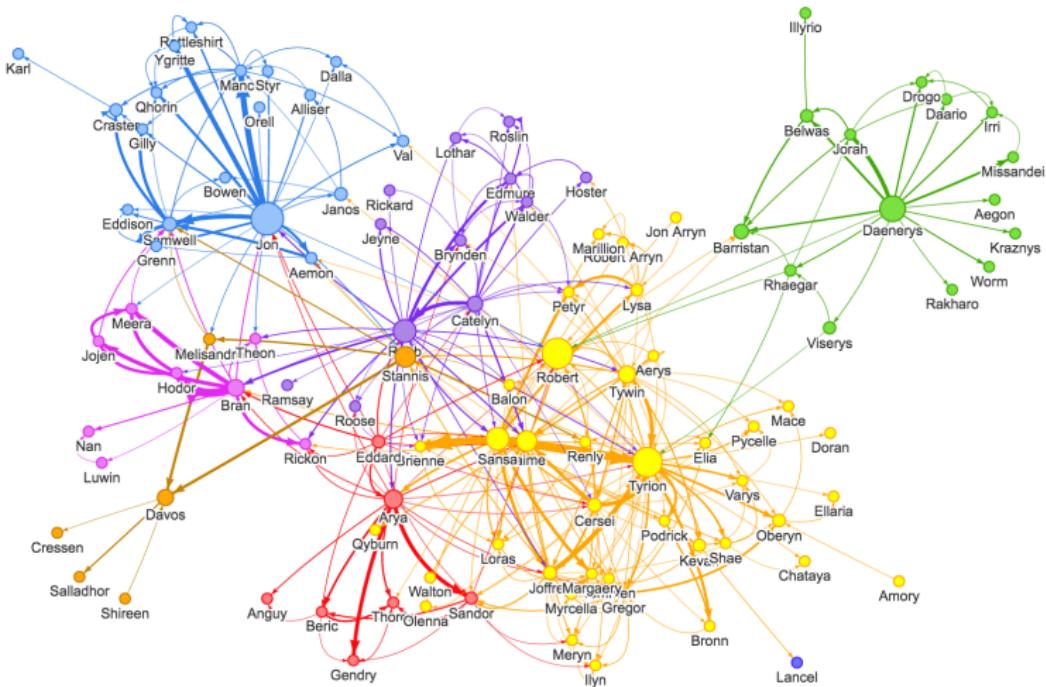
$G = \text{Friendship Network}$ ,  $V = \text{persons}$ ,  $E = \text{friendship}$

$G = \text{Airport Network}$ ,  $V = \text{airports}$ ,  $E = \text{non-stop flight}$

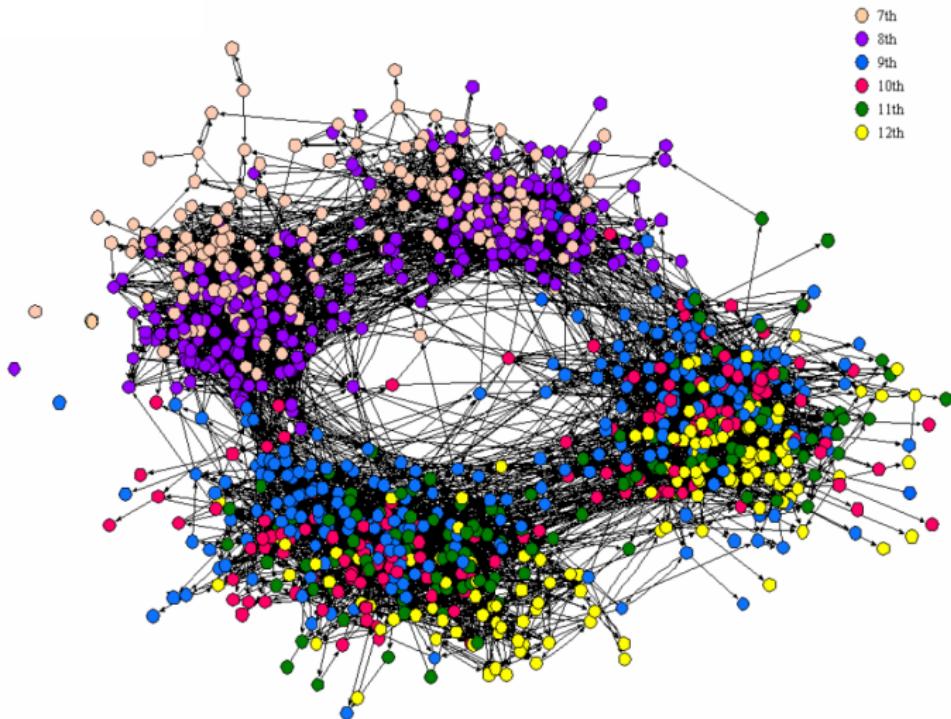
# Worldwide Airports (2017)



# 'Friendship' Network (2015)

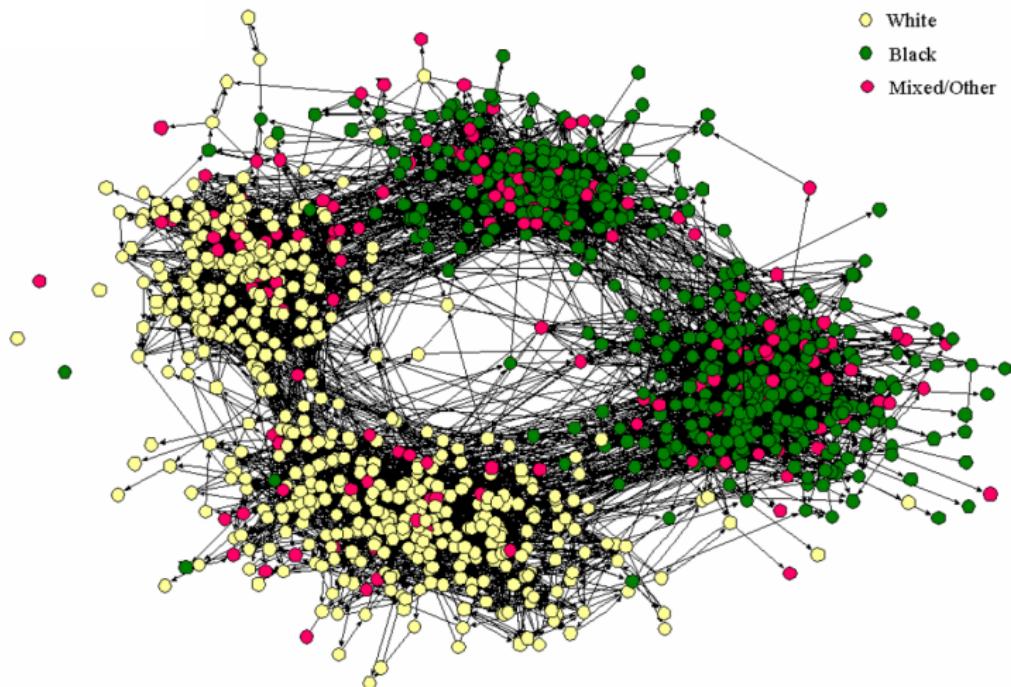


# High School Friendship (2001)



from Moody (2001) Race, School Integration and Friendship

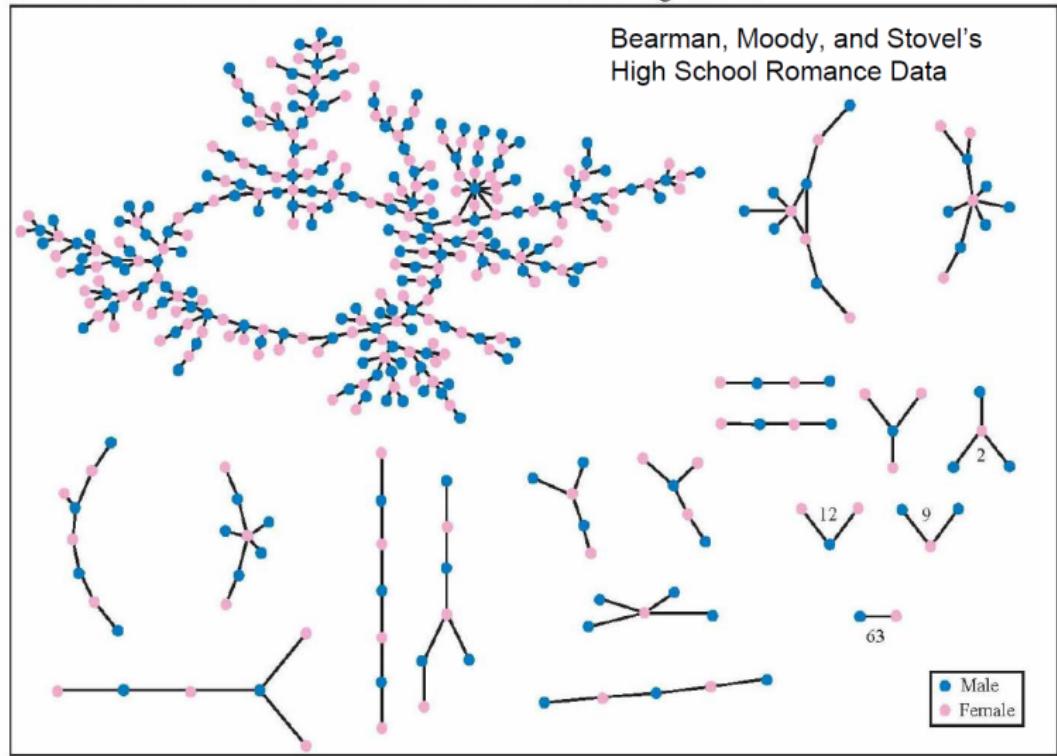
# High School Friendship (2001)



from Moody (2001) Race, School Integration and Friendship

# Romantic and Sexual Relationships

The Structure of Romantic and Sexual Relations at "Jefferson High School"



# Cyberisk in Network Models

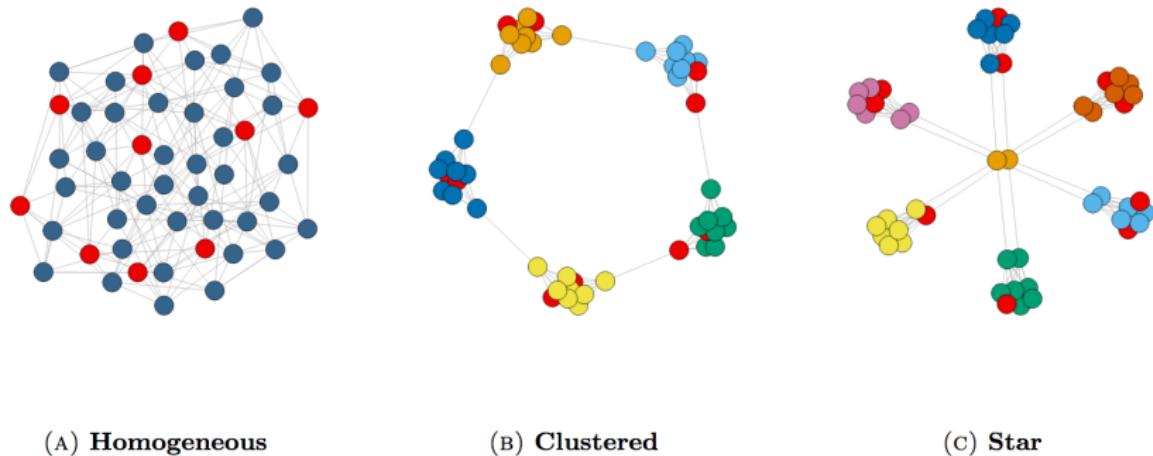
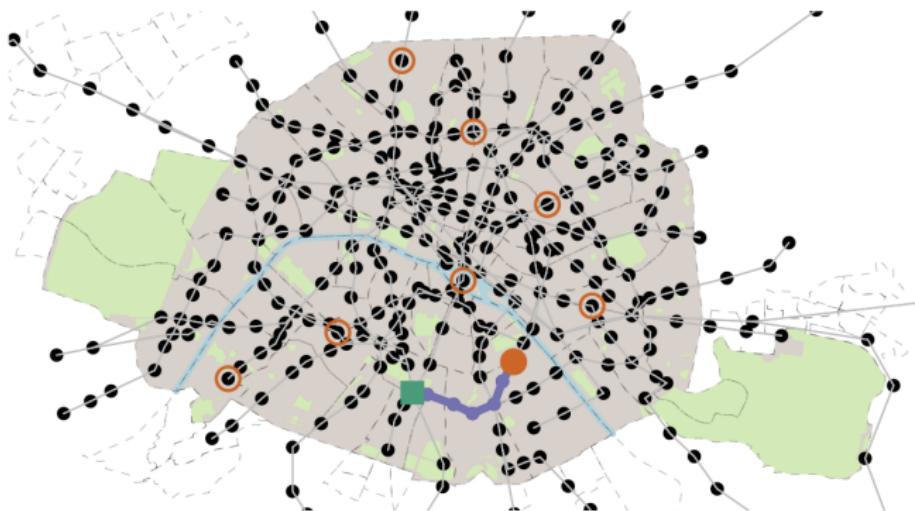


FIGURE 6. Infection scenario: The red nodes (3,5,13,15,23,25,33,35,43,45) are initially infected.

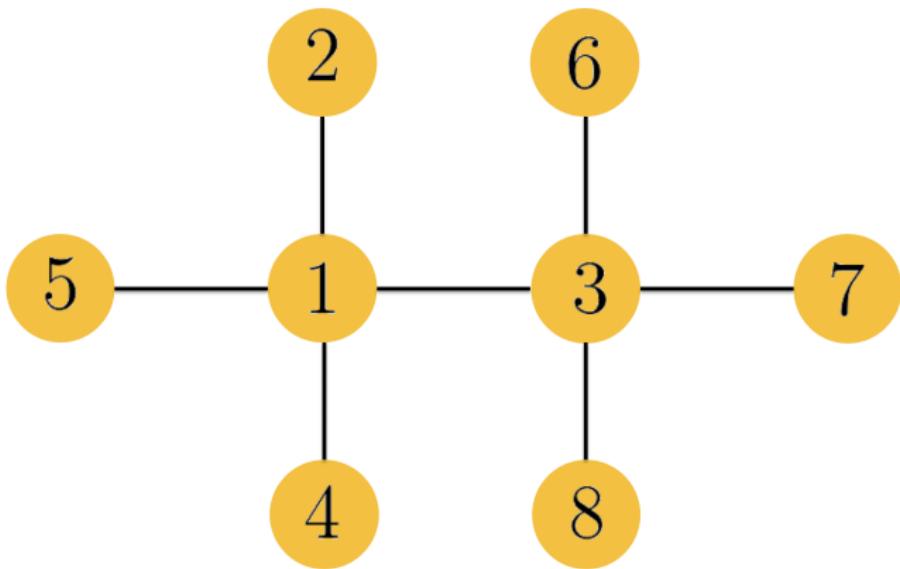
Fahrenwaldt *et al.* (2017, Pricing of cyber insurance contracts)

## Shortest Path

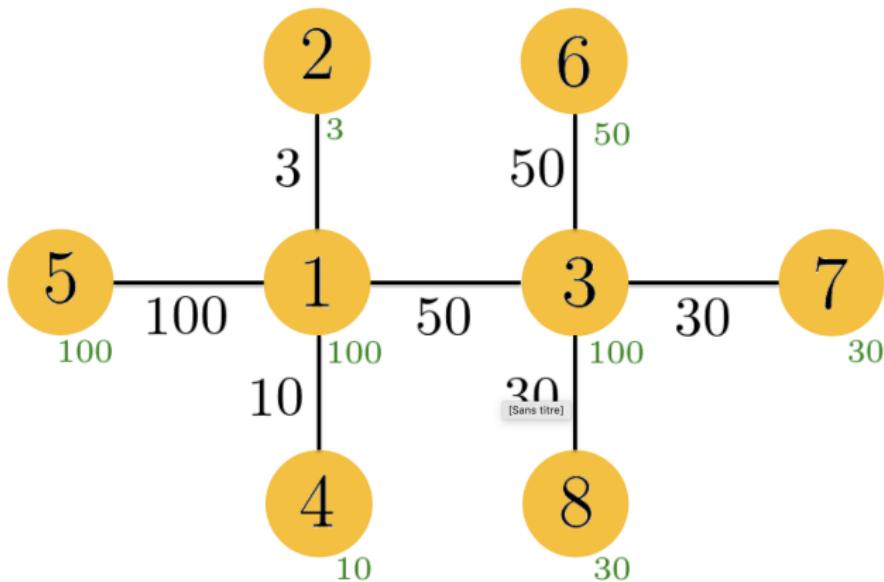


Charpentier et al. (2019, Optimal transport on large networks: a practitioner guide)

# Networks & Insurance (P2P)



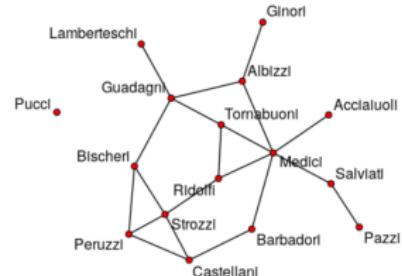
## Networks & Insurance (P2P)



used by Inspeer, see <https://wwwthe-digital-insurer.com>

# Networks with R

```
1 > library(network)
2 > data(flo)
3 > nflo<-network(flo,directed=FALSE)
4 > nflo
5 Network attributes:
6   vertices = 16
7   directed = FALSE
8   total edges= 20
9 > plot(nflo,displaylabels=TRUE,
       boxed.labels=FALSE)
```



Florentine Wedding Data.  $n_V = 16$  vertices (families),  $n_E = 20$  edges (weddings), from Padgett (1994, *Marriage and Elite Structure in Renaissance Florence, 1282-1500*)

# Network Representation

Let  $V = \{1, \dots, n_V\}$  denote either **nodes**, or **vertices** ( $n_V$  is the **order**). Let  $E \in \{0, 1\}^{n_V \times n_V}$  represents the **relationships**, through an adjacency matrix  $A$ ,  $A_{i,j} = 1$  indicates a **link** - or **edge** - between  $i$  and  $j$ , or a collection of links  $\{e_1, \dots, e_{n_E}\}$ . Let  $n_E = |E|$  denote the number of edges, called **size**.

The **degree**  $d(\cdot)$  of a vertex  $v$  is its number of incident edges.

A **network** (or a **graph**) is a pair  $G = (V, E)$

If  $V_1 \subset V_2$  and  $E_1 \subset E_2$ , then  $(V_1, E_1)$  is a **subgraph** of  $(V_2, E_2)$ .

Two vertices  $u, v \in V$  are said to be **adjacent** (or **connected**) if they are joined by an edge in  $E$ .

## Network Representation

One might consider some **undirected network**: consider e.g. **adjacency matrix  $A$**

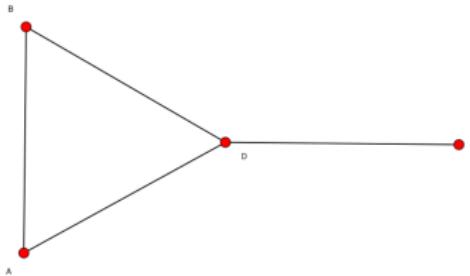
$$A = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

with  $n_V = 4$  vertices, and  $n_E = 4$  edges (4 1's in the upper corner of the matrix).

There are no self-loops, i.e.  $A_{i,i} = 0$ .

## Network Representation

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



$A_{ij} = 1$  if and only if  $i$  and  $j$  are linked. The matrix is symmetric ( $A_{ij} = A_{ji}$ ), the network is **undirected**.

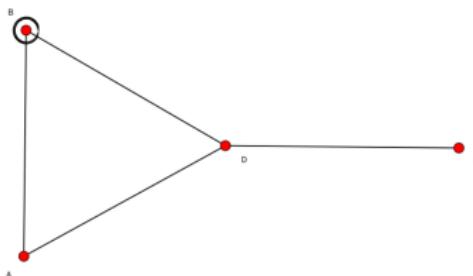
Links are  $E = \{1 - 2, 1 - 4, 2 - 4, 3 - 4\}$  - no need to mention  $\{2 - 1, 4 - 1, 4 - 2, 4 - 3\}$  since undirected. Hence,  $n_E = 4$ .  
Further,  $d(1) = 2$ ,  $d(2) = 2$ ,  $d(3) = 1$  and  $d(4) = 3$ .

# Network Representation

Row  $i$  contains list of vertices connected to vertice  $i$ .

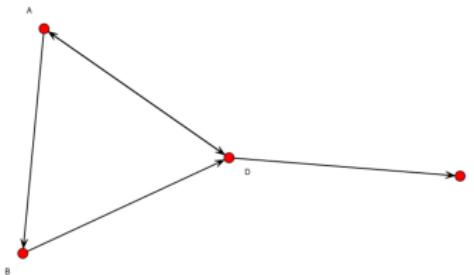
$$d(i) = \sum_{j=1}^{n_V} A_{i,j} = \mathbf{A}_{i,\cdot}^T \mathbf{1}.$$

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



# Network Representation

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



Here the network is **directed** (also called **digraph**).

Links are  $E = \{(1, 2), (1, 4), (2, 4), (4, 1), (4, 3)\}$  (the 5 arrows).

or  $E = \{(1 \rightarrow 2), (1 \leftrightarrow 4), (2 \rightarrow 4), (3 \leftarrow 4)\}$

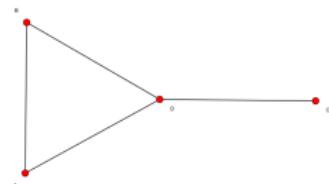
For edge  $(1, 2)$ , 1 is the **source** node and 2 is the **terminal** node.

Vertices have **in-degrees** and **out-degrees**.

## Network Representation

For the undirected graph, one can also consider the  $n_V \times n_E$  incidence matrix  $\mathbf{T}$

$$\mathbf{T} = \begin{matrix} & \begin{matrix} (ab) & (ad) & (bd) & (cd) \end{matrix} \\ \begin{matrix} a \\ b \\ c \\ d \end{matrix} & \left( \begin{matrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \end{matrix} \right) \end{matrix}$$



Observe that  $d(v) = \sum_{e=1}^{n_E} T_{v,e}$ , i.e.  $\mathbf{d} = \mathbf{T}\mathbf{1}$ .

Further  $\mathbf{T}\mathbf{T}^T = \mathbf{D} + \mathbf{A}$  where  $\mathbf{D} = \text{diag}[\mathbf{d}]$

## Network Representation

Finally, the  $n_V \times n_V$  matrix  $\mathbf{L} = \mathbf{D} - \mathbf{A}$  is called the **Laplacian** of the graph. One can prove that for any  $\mathbf{x} \in \mathbb{R}^{n_V}$ ,

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (x_i - x_j)^2$$

The normalized Laplacian is  $\tilde{\mathbf{L}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$ . Note that eigenvalues of  $\tilde{\mathbf{L}}$  lie in interval  $[0, 2]$ .

## Network Representation

The goal is to embed a combinatorial object, our network,  $(V, E)$  into a two-dimensional Euclidean space.

Clearly **not unique**...

The visualization of network  $G = (V, E)$  is a mapping  
 $\varphi : (V, E) \rightarrow \mathbb{R}^2$ . No geometry in  $G$

**Multidimensional scaling** (MDS) is commonly used for visualization. Given vertex distances  $D = [D_{i,j}]$  (shortest path, e.g.), we want to find

$$\vec{z}_i = (x_i, y_i), \text{ for } i \in V, \text{ so that } \|\vec{z}_i - \vec{z}_j\| \sim D_{i,j}$$

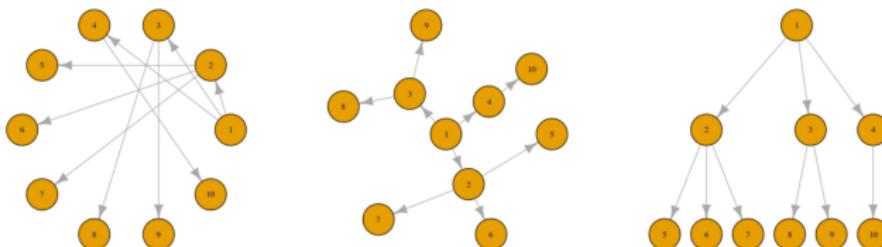
Consider

$$\min_{\vec{z}_1, \dots, \vec{z}_{n_V}} \left\{ (D_{i,j} - \|\mathbf{z}_i - \mathbf{z}_j\|)^2 \right\}$$

Possible to add some constraints, e.g. centralities

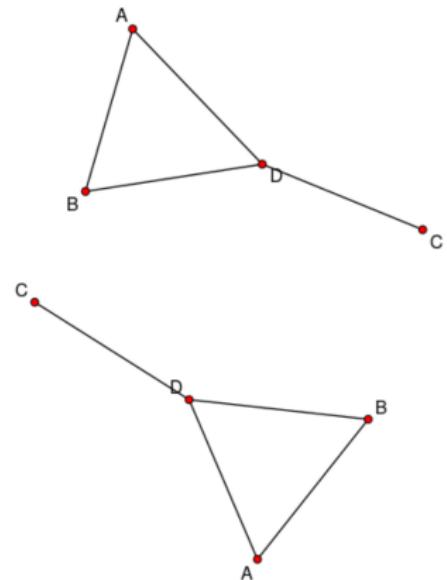
# Network Representation

- convention : straight line segments
- aesthetics : minimal edge crossing
- aesthetics : relative placement of vertices, subgraphs, etc.



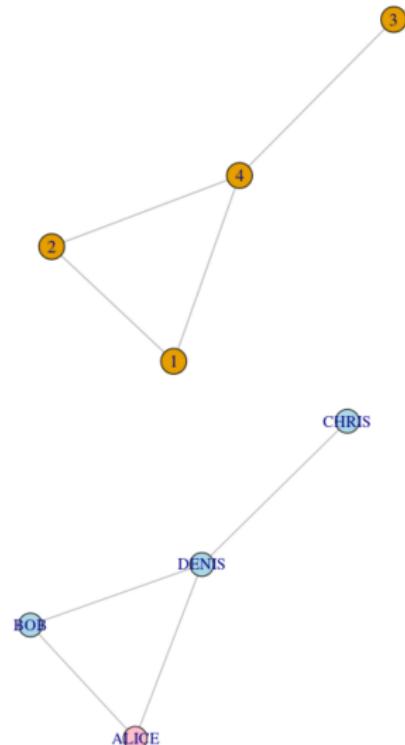
# Networks with R

```
1 > library(network)
2 > m = matrix(0,4,4)
3 > colnames(m)=rownames(m)=LETTERS
   [1:4]
4 > m[1,2]=m[1,4]=m[2,1]=m[2,4]=m
   [3,4]=m[4,1]=m[4,2]=m[4,3]=1
5 > m
6   A B C D
7 A 0 1 0 1
8 B 1 0 0 1
9 C 0 0 0 1
10 D 1 1 1 0
11 > nm=network(m,directed=FALSE)
12 > plot(nm,displaylabels=TRUE,boxed.
   labels=FALSE)
```



# Networks with R

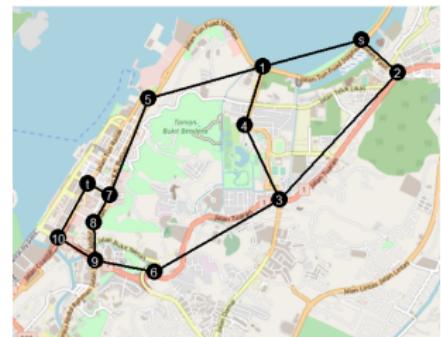
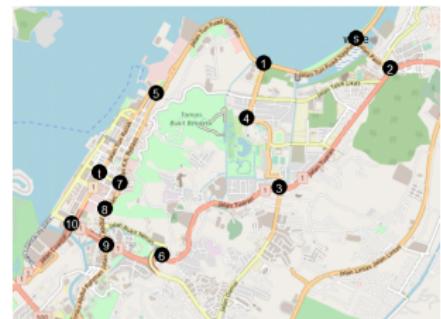
```
1 > library(igraph)
2 > g <- graph.formula
   (1-2,1-4,2-4,3-4)
3 > plot(g)
4 > V(g)
5 + 4/4 vertices, named, from 3b97949
  :
6 [1] 1 2 4 3
7 > V(g)$label=c("ALICE","BOB","DENIS",
   "CHRIS")
8 > V(g)$color=c("pink",rep("light
   blue",3))
9 > E(g)
10 + 4/4 edges from 3b97949 (vertex
    names):
11 [1] 1--2 1--4 2--4 4--3
12 > plot(g)
```



# Spatial Data & Graphs

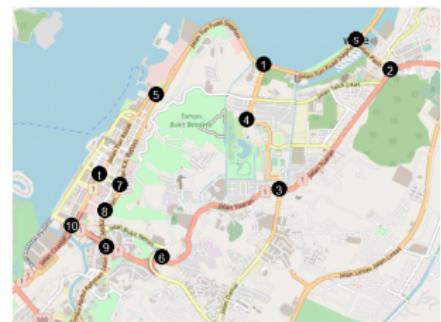
From Abdullah & Kien Hua (2017, [The Application of the Shortest Path and Maximum Flow with Bottleneck in Traffic Flow](#))

```
1 location <- 'http://www.jistm.com/  
PDF/JISTM-2017-04-06-02.pdf'  
2 m <- matrix(c(0, 5.995910,  
116.105520, 1, 5.992737,  
116.093718, 2, 5.992066,  
116.109883, 3, 5.976947,  
116.095760, 4, 5.985766,  
116.091580, 5, 5.988940,  
116.080112, 6, 5.968318,  
116.080764, 7, 5.977454,  
116.075460, 8, 5.974226,  
116.073604, 9, 5.969651,  
116.073753, 10, 5.972341,  
116.069270, 11, 5.978818,  
116.072880), 3, 12)
```

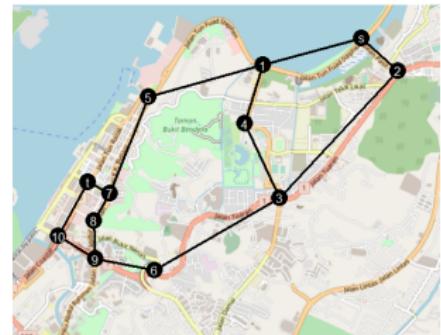


# Spatial Data & Graphs

```
1 library(OpenStreetMap)
2 map <- openproj(openmap(c(lat=
6.000, lon= 116.06), c(lat=
5.960, lon= 116.12)))
3 plot(map)
4 points(t(m[3:2,]), col="black", pch
=19, cex=3)
```



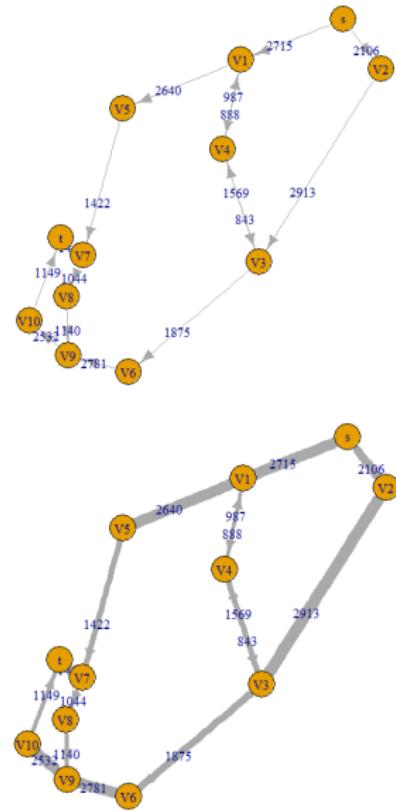
```
1 library(tabulizer)
2 out <- extract_tables(location)
3 B1 = as.data.frame(out[[2]])
4 B2 = as.data.frame(out[[3]])
5 E = data.frame(from=B1[3:20,"V3"],
to=B1[3:20,"V4"])
6 E = E[-c(6,8),]
7 capacity = as.character(B2$V3[-1])
8 E$capacity = as.numeric(capacity)
```



# Spatial Data & Graphs

We can plot only the network  
(without background map)

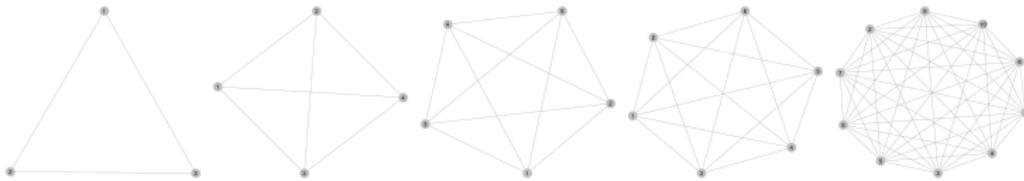
```
1 library(igraph)
2 g=graph_from_data_frame(E)
3 E(g)$label=E$capacity
4 plot(g, layout=as.matrix(B[,c("x",
5 y")]))
6 plot(g, layout=as.matrix(B[,c("x",
7 y")]), edge.width=E$capacity
8 /200)
```



## Special Graphs

A **complete graph** of size  $n$  has  $n$  vertices and  $\frac{n(n - 1)}{2}$  edges.

$$A_{i,j} = \mathbf{1}_{i \neq j}.$$



A **regular graph** is a graph in which every vertex has the same degree.



## Random Graphs

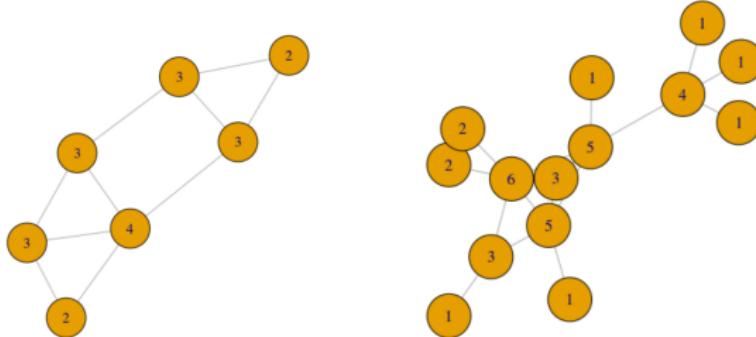
The **degree** of a vertex  $v$  is  $d(v)$ , the number of vertices in  $V$  incident to  $v$  (i.e. the number of neighbors of  $v$ )

A network is said to be **sparse** if

$$n_E \ll \frac{n_V(n_V - 1)}{2} \text{ or } \bar{d} = \frac{1}{n_V} \sum_{v=1}^{n_V} d_v \ll n_V - 1$$

Let  $f_\delta$  denote the proportion of vertices with degree  $\delta$ , and  $\mathbf{f} = (f_\delta)$  the degree distribution, which is a summary of local connectivity across the graph.

# Random Graphs

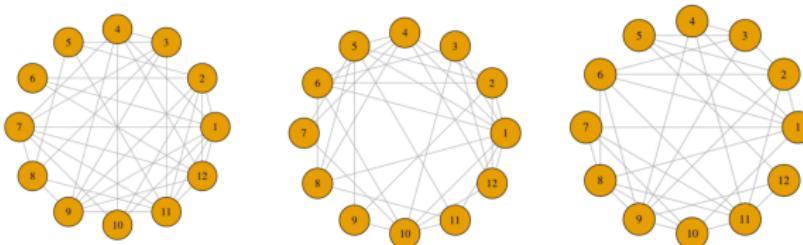


High-degree vertices are likely to be influential, central, prominent.

# Networks with R

The Erdős-Renyi random graph model  $G_{n,p}$  is an undirected graph with  $n$  vertices, such that edge  $(v, v')$  is present with probability  $p$ , independent of other edges.

```
1 library(igraph)
2 E = erdos.renyi.game(n=12, p.or.m=.5, type="gnp")
3 plot(E, layout=layout.circle)
4 # hist(igraph::degree(E))
```



## Random Graphs

Then  $d(v)$  is a binomial distribution  $\mathcal{B}(n_V - 1, p)$ , i.e.

$$f_\delta = \mathbb{P}[d(v) = \delta] = \binom{n-1}{\delta} p^\delta (1-p)^{n_V-1-\delta}$$

For a large network ( $n_V \rightarrow \infty$ ),  $d(v) \sim \mathcal{N}(n_V p, n_V p(1-p))$  from the law of large numbers.

For a large network with  $p \sim \lambda/n_V$ ,  $d(v) \sim \mathcal{P}(\lambda)$  from the law of small numbers.

## Random Graphs: Power Laws

Scale-free network have degree distribution with power-law tail

A scale free function  $f(\cdot)$  satisfies  $f(ax) = bf(x)$ ,  $\forall x$ , for some  $a, b$

**Ex** Power-law functions  $f(x) = x^{-\alpha}$  are scale-free

$$f(ax) = [ax]^{-\alpha} = a^{-\alpha}f(x) = bf(x), \text{ where } b = a^{-\alpha}$$

$\log f_\delta \sim c - \alpha \log \delta$  for some constant  $c$

Power-law exponent (negative slope) is typically  $\alpha \in [2, 3]$

The normalized power-law degree distribution is

$$f_\delta = \frac{\alpha - 1}{\delta_0} \left( \frac{\delta}{\delta_0} \right)^{-\alpha}, \text{ for } \delta \geq \delta_0.$$

More convenient to assume (for computation) that  $\delta \in \mathbb{R}$  (instead of  $\mathbb{N}$ )

## Random Graphs: Power Laws

For instance, the probability that a random node has degree exceeding 100 is

$$\mathbb{P}[f_\delta > 100] = \int_{100}^{\infty} \frac{\alpha - 1}{\delta_0} \left(\frac{x}{\delta_0}\right)^{-\alpha} dx = \left(\frac{100}{\delta_0}\right)^{1-\alpha}$$

(which is also a power function, with index  $\alpha - 1$ ).

Further

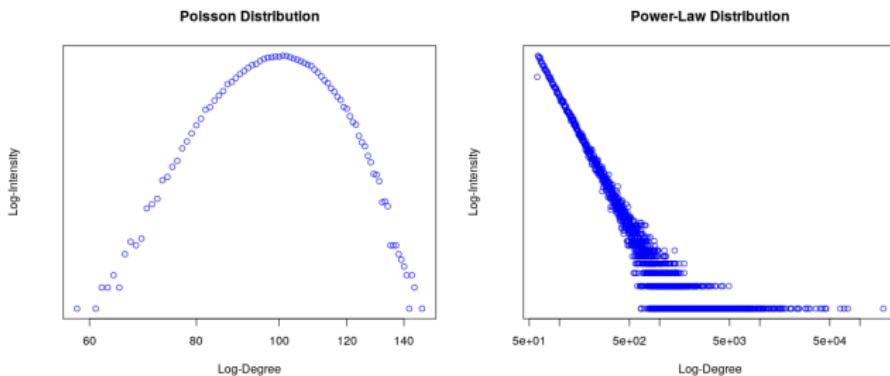
$$\mathbb{E}[f_\delta] = \frac{\alpha - 1}{\alpha - 2} \delta_0$$

but if  $\alpha \in [2, 3]$ ,  $\text{Var}[f_\delta] = \infty$ .

# Networks with R

Use log-log Pareto plots to visualize  $f_\delta$ , either on  $\log f_\delta$ ,

```
1 > dd.yeast <- table(degree(g))/n
2 > ind <- (dd.yeast != 0)
3 > d <- as.numeric(names(ind))
4 > plot(d[ind], dd.yeast[ind], log="xy")
```



# Networks with R

... or on  $\log \bar{F}_\delta$  (survival cumulative distribution)

```
1 > dd.yeast <- table(degree(g))/n
2 > ind <- (dd.yeast != 0)
3 > d <- as.numeric(names(ind))
4 > plot(d[ind], 1-cumsum(dd.yeast[ind]), log="xy")
```

(classical in extreme value theory: log-log Pareto plot)

## Random Graphs: Power Laws

Natural to consider the linear least-squares (LS) estimator of  $c$  and  $\alpha$ ,

$$\min \left\{ \sum_{\delta} (\log \bar{F}_{\delta} - c + [\alpha - 1] \log \delta)^2 \right\}$$

Popular, but extremely noisy, biased (log transformations), and valid only when  $\delta > \delta_0$

Here

$$f_{\delta} = \frac{\alpha - 1}{\delta_0} \left( \frac{\delta}{\delta_0} \right)^{-\alpha} \text{ for } \delta \geq \delta_0,$$

so that the log-likelihood function is (up to constants independent of  $\alpha$ )

$$\log \mathcal{L}(\alpha) = \sum_{v=1}^{n_V} \log f_{\delta_v} \propto n_V \log[\alpha - 1] - \alpha \sum_{v=1}^{n_V} \log \left( \frac{\delta_v}{\delta_0} \right)$$

## Random Graphs: Power Laws

and the maximum likelihood estimator is then

$$\hat{\alpha} = 1 + \left( \frac{1}{n_V} \sum_{v=1}^{n_V} \log \left( \frac{\delta_v}{\delta_0} \right) \right)^{-1}$$

also called **Hill estimator**

Usually, we consider the  $k$  largest values (so-called **tail**)

$$\hat{\alpha}_k = 1 + \left( \frac{1}{k} \sum_{i=0}^{k-1} \log \left( \frac{\delta_{n_V-i:n_V}}{\delta_{n_V-k:n_V}} \right) \right)^{-1}$$

Hill's plot is the graph of  $\hat{\alpha}_k$  as a function of  $k$

in standard litterature on **extreme value** the so called  $\alpha$  is  $\alpha - 1$  from the network literature...

## Preferential attachment model

Here also it is possible to derive a stochastic network representation.

Classical model for popularity, see Yule (1925, *A Mathematical Theory of Evolution*) or Merton (1968, *The Matthew effect in science*).

- with probability  $p$ ,  $v'$  connects to  $v$  randomly (uniformly),
- with probability  $1 - p$ ,  $v'$  connects to  $v$  with **probability proportional to the degree of  $v$**

Preferential attachment model leads to **rich-gets-richer** dynamics

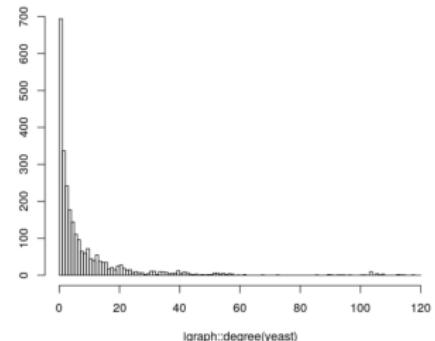
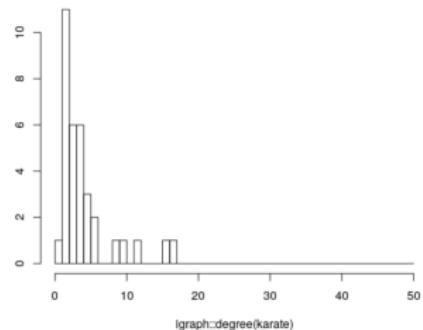
Degrees have a power-law distribution with tail exponent

$$\alpha = 1 + \frac{1}{1-p}$$

$$f_\delta \sim \frac{1}{p} \left( \left( \frac{1-p}{p} \right) \delta + 1 \right)^{-\left(1+\frac{1}{1-p}\right)}$$

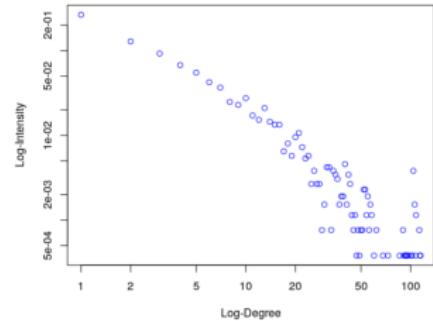
# Networks with R

```
1 data(karate)
2 hist(igraph::degree(karate),breaks=
      seq(0,50))
3 data(yeast)
4 hist(igraph::degree(yeast),breaks=
      seq(0,120))
```



# Networks with R

```
1 dd.yeast = degree.distribution(  
  yeast)  
2 d = 1:max(d.yeast)-1  
3 ind = (dd.yeast != 0)  
4 plot(d[ind], dd.yeast[ind], log="xy"  
  , col="blue", xlab=c("Log-  
  Degree"), ylab=c("Log-Intensity  
  "))
```



## Paths in Networks

A **walk** of length  $k$  from vertex  $v_0$  to vertex  $v_k$  is a sub-graph  $(V_p, E_p)$  with  $V_p = \{v_0, v_1, \dots, v_k\}$  and  $E_p = \{(v_0, v_1), (v_1, v_2), \dots, (v_{k-1}, v_k)\}$ .

A **trail** is a walk without repeated edges

A **path** is a walk without repeated vertices, and therefore with different edges.

If  $v_0 = v_k$ , a walk is said to be closed. A **cycle** is a closed walk with different vertices (except  $v_0 = v_k$ ), also called **circuit**.

A walk from  $v_0$  to  $v_k$  contains a path from  $v_0$  to  $v_k$  (remove subcycles).

If there is a walk from  $v$  to  $v'$ , then  $v$  and  $v'$  are **connected** - or  $v'$  is **reachable** from  $v$ . A graph is **connected** if every vertex is reachable from every other.

## Paths in Networks

Consider a network with 3 nodes and adjacency matrix



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

```
1 > A=matrix(c(0,1,1,1,0,0,1,0,0),3,3)
2 > A%*%A
3   [,1]  [,2]  [,3]
4 [1,]    2      0      0
5 [2,]    0      1      1
6 [3,]    0      1      1
```

## Paths in Networks

**Proposition**  $[\mathbf{A}^k]_{i,j}$  is the number of walks of length  $k$  from  $i$  to  $j$ .  
Indeed, there are 2 paths of length 2 from 1 to 1 :  $(1,2)+(2,1)$  or  $(1,3)+(3,1)$

One can prove that

$$\mathbf{A}^{2k} = \begin{pmatrix} 2^k & 0 & 0 \\ 0 & 2^{k-1} & 2^{k-1} \\ 0 & 2^{k-1} & 2^{k-1} \end{pmatrix} \text{ and } \mathbf{A}^{2k+1} = \begin{pmatrix} 0 & 2^k & 2^k \\ 2^k & 0 & 0 \\ 2^k & 0 & 0 \end{pmatrix}$$

**Exploration Algorithm:** find the set of all vertices that can be reached by a walk from  $v \in V$ , denoted  $\mathcal{C}(v)$ .

Because of cycle properties, repeat exploration for vertices in  $V \subset \mathcal{C}(v)$ .

# Connectivity in Networks

How to test for connectivity in a graph?

Idea : use adjacency list from a starting vertice  $s$  to explore

- Set  $L = M = \{s\}$  and repeat while there are still nodes to explore ( $L \neq \emptyset$ )
  - Pick  $u \in L$ 
    - If there is an edge  $(u, v) \in E$  with  $v \in M$  then select one , and set  $L = L \cup \{v\}$  and  $M = M \cup \{v\}$
    - Otherwise prune, i.e.  $L = L \setminus \{v\}$

Exploration is of order  $2nV$ , each node is added and removed once.

# Connectivity in Networks

$L$	$M$
{2}	2
{2, 1}	1
{2, 1, 5}	5
{2, 1, 5, 6}	6
{1, 5, 6}	
{1, 5, 6, 4}	4
{5, 6, 4}	
{5, 4}	
{5, 4, 3}	3
{5, 3}	
{5, 3, 7}	7
{5, 3}	
{3}	
{3, 8}	8
{3}	
{}	

# Connectivity in Networks (BFS)

Breadth-first search (BFS):  $u$  should be the first element of  $L$

$L$	$M$
$\{2\}$	2
$\{2, 1\}$	1
$\{2, 1, 5\}$	5
$\{1, 5\}$	
$\{1, 5, 4\}$	4
$\{1, 5, 4, 6\}$	6
$\{5, 4, 6\}$	
$\{4, 6\}$	
$\{4, 6, 3\}$	3
$\{6, 3\}$	
$\{3\}$	
$\{3, 7\}$	7
$\{3, 7, 8\}$	8
$\{7, 8\}$	
$\{8\}$	
$\{\}$	

# Connectivity in Networks (DFS)

Depth-first search (DFS):  $u$  should be the last element of  $L$

$L$	$M$
{2}	2
{2, 1}	1
{2, 1, 4}	4
{2, 1, 4, 3}	3
{2, 1, 4, 3, 7}	7
{2, 1, 4, 3}	
{2, 1, 4, 3, 8}	8
{2, 1, 4, 3}	
{2, 1, 4}	3
{2, 1, 4, 6}	6
{2, 1, 4, 6, 5}	5
{2, 1, 4, 6}	
{2, 1, 4}	
{2, 1}	
{2}	
{}	

## Connectivity in Networks

Another important concept is related to **resilience** of networks to the removal of some vertices.

As vertices are removed, shortest paths distances increase (can even be infinite when the network is disconnected).

Important in epidemiology and vaccination models.

## Connectivity in Networks

The **shortest path** between two vertices (or nodes) in a graph is such that the number of its constituent edges is minimized.

This **shortest path** is often referred as the **geodesic distance**

The longest shortest path in a graph is the **diameter** of the graph  
more generally, the sum of the weights of its constituent edges is minimized.

## Connectivity in Networks

(initially designed for a digraph, with possible costs)

Start with  $\mathcal{S} = \{v\}$ ,  $\delta(v, v') = 1$  if  $A_{v,v'} = 1$ ,  $\infty$  otherwise.

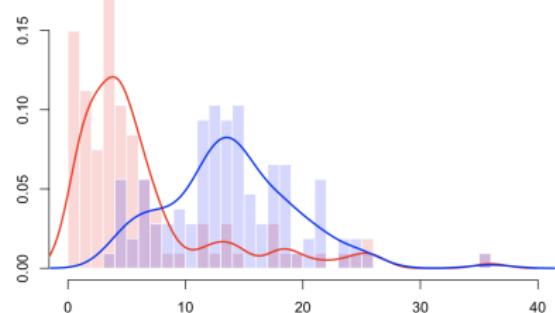
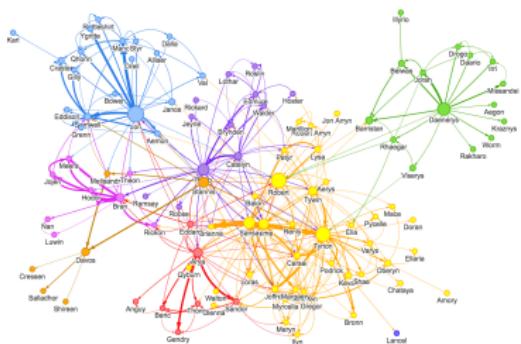
Until  $\mathcal{S} = V$  : if there is  $v' \notin \mathcal{S}$  such that  $A_{u,v'} = 1$ ,  $\forall u \in \mathcal{S}$ , then  $\mathcal{S} = \mathcal{S} \cup \{v'\}$  and

$$\delta(v') = \min\{\delta(v'), \delta(u) + 1\}$$

See [Dijkstra's algorithm](#), or Floyd–Warshall algorithm

# Friendship Paradox

*“The average degree of a friend is strictly greater than the average degree of a random node”, see Feld (1991) and Zuckerman & Jost (2001).* See Game of Thrones network



## Friendship Paradox

The average number of friends of a random person in the graph is

$$\mu = \frac{1}{n_V} \sum_{v \in V} d(v) = \frac{2n_E}{n_V} = \frac{1}{n_V} \|\mathbf{d}\|_1$$

where  $\mathbf{d}$  is the vector of  $d(v)$ 's (or  $\mathbf{d} = \mathbf{A}\mathbf{1}$ ).

The average number of friends that a typical friend has is

$$\frac{1}{n_V} \sum_{v \in V} \left( \frac{1}{d(v)} \sum_{v' : (v, v') \in E} d(v') \right)$$

Observe that each node  $v$  appears as a neighbour  $d(v)$  times and each time it contributes to  $d(v)$  neighbours. Since we simply count how many neighbours neighbours have, this average can also be written

$$\frac{1}{\|\mathbf{d}\|_1} \mathbf{d}^\top \mathbf{d}$$

## Friendship Paradox

From Cauchy-Scharwz, for any  $\mathbf{u}$  and  $\mathbf{v}$ ,  $\mathbf{u}^\top \mathbf{v} \leq \|\mathbf{u}\|_2 \cdot \|\mathbf{v}\|_2$ . If  $\mathbf{u} = \mathbf{1}$  and  $\mathbf{v} = \mathbf{d}$ , we can write

$$\mathbf{1}^\top \mathbf{d} = \|d\|_1 \leq \|\mathbf{1}\|_2 \cdot \|\mathbf{d}\|_2 = \sqrt{n_V} \cdot \sqrt{\mathbf{d}^\top \mathbf{d}}$$

After squaring and rearranging, we get

$$\frac{1}{\|\mathbf{d}\|_1} \mathbf{d}^\top \mathbf{d} \geq \frac{1}{n_V} \|\mathbf{d}\|_1$$

Alternatively, for a random variable  $D$  denoting the number of friends of a random individual, and  $D^2$  is seen as number of friends of friends. Since  $\text{Var}[D] = \mathbb{E}[D^2] - \mathbb{E}[D]^2$  we can write

$$\frac{\mathbb{E}[D^2]}{\mathbb{E}[D]} = \mathbb{E}[D] + \frac{\text{Var}[D]}{\mathbb{E}[D]} \geq \mathbb{E}[D]$$

because  $\text{Var}[D] \geq 0$ .

## Friendship Paradox

*“You apply for a loan and your would-be lender somehow examines the credit ratings of your Facebook friends. If the average credit rating of these members is at least a minimum credit score, the lender continues to process the loan application. Otherwise, the loan application is rejected,” the patent states.”* Bhattacharya (2015, Facebook patent: Your friends could help you get a loan - or not)

Consider some positive variable  $y$  observed at each node  $v$ . It can be the credit score of friends, or activity of people on twitter, as in Hodas *et al.* (2013, Friendship paradox redux: Your friends are more interesting than you)

## Friendship Paradox

Using the extension of König-Huyghens,

$\text{Cov}[D, Y] = \mathbb{E}[DY] - \mathbb{E}[D]\mathbb{E}[Y]$ , write  $\mathbf{d}^\top \mathbf{y}$  as

$$\mathbf{d}^\top \mathbf{y} = \frac{\|\mathbf{d}\|_1 \|\mathbf{y}\|_1}{n_V} + n_V \text{Cov}[D, Y]$$

or

$$\frac{1}{\|\mathbf{d}\|_1} \mathbf{d}^\top \mathbf{y} = \frac{1}{n_V} \|\mathbf{y}\|_1 + \frac{n_V}{\|\mathbf{d}\|_1} \text{Cov}[D, Y]$$

This can be written also

$$\frac{\mathbb{E}[DY]}{\mathbb{E}[D]} = \mathbb{E}[Y] + \frac{\text{Cov}[D, Y]}{\mathbb{E}[D]}$$

Hence,

$$\frac{\mathbb{E}[DY]}{\mathbb{E}[D]} = \mathbb{E}[Y] \text{ if } X \perp Y \text{ and } \frac{\mathbb{E}[DY]}{\mathbb{E}[D]} > \mathbb{E}[Y] \text{ if } \text{Cov}[D, Y] > 0.$$

## Friendship Paradox

This concept is related to the *generalized friendship paradox with respect to quantity  $y$*  from Em & Jo (2014, [Generalized friendship paradox in complex networks: The case of scientific collaboration](#)) which states that it will arise if  $y$  is nonnegatively correlated with degree  $d$ , or Charpentier & Galichon (2019).

## Networks & Centrality

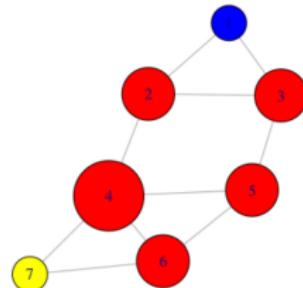
The goal is to understand who is important, based on their network position.

*“There is certainly no unanimity on exactly what centrality is or on its conceptual foundations, and there is little agreement on the proper procedure for its measurement”*, Freeman (1979, *Centrality in Social Networks*)

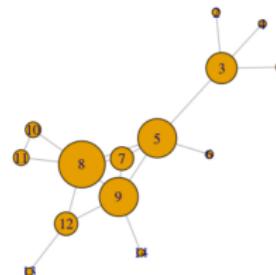
Simple question such as *which vertex is central* can be replaced by more interesting ones, such as *what percentage of vertices is crucial to the network connectivity ?*

He or she who has many friends is most important, divided by the maximum possible  $n_V$  : for a vertice  $v$ ,  $C_D(v) = \frac{d(v)}{n_V}$

# Networks with R

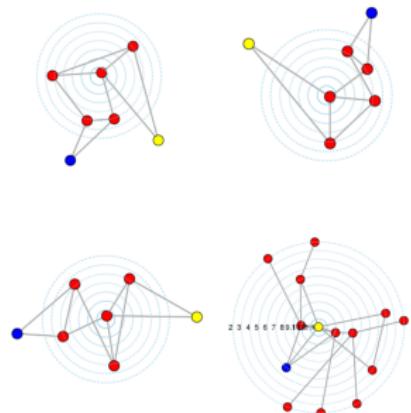


```
1 plot(g, vertex.size=degree(ng))
```



# Networks with R

```
1 A = get.adjacency(g1, sparse=FALSE)
2 ng = network::as.network.matrix(A)
3 sna::gplot.target(ng, degree(ng),
                  circ.col="skyblue", vertex.col=c
                  ("blue", rep("red", 5), "yellow
                  "), edge.col="darkgray")
4 degree(ng)
5 [1] 4 6 6 8 6 6 4
```



## Networks & Centrality

How much variation is there in the centrality scores among the nodes?

One can define a centralization index, using any dispersion measure (variance, Gini index).

Freeman's centralization is based to the distance to the maximum

$$C_D = \frac{1}{(n_V - 1)(n_V - 2)} \sum_{v \in V} [d(v^*) - d(v)]$$

where  $v^* = \operatorname{argmax}\{d(v)\}$  (supposed to be unique).

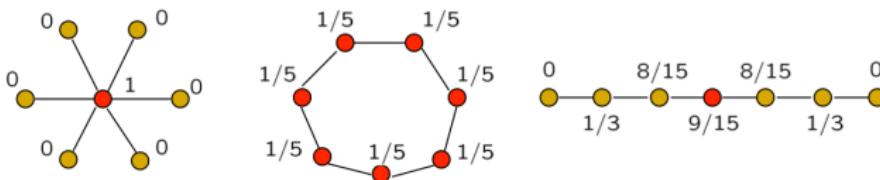
## Networks & Centrality

How many pairs of individuals would have to go through you in order to reach one another in the minimum number of hops?

The betweenness centrality of a node  $v$  is

$$C_B(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

where  $\sigma_{st}$  is the total number of shortest paths from node  $s$  to  $t$ , and  $\sigma_{st}(v)$  is the number of those paths that pass through  $v$ .



## Networks & Centrality

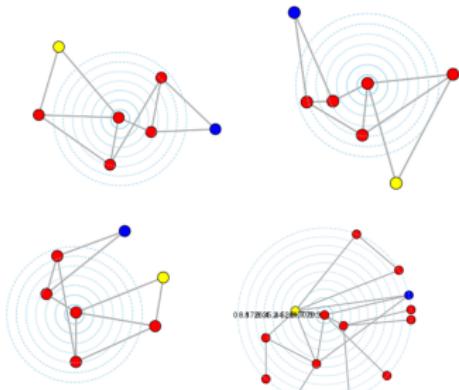
One possible normalization is  $C_{B'}(v) = \frac{2C_B(v)}{(n_V - 1)(n_V - 2)}$   
(number of pairs of vertices excluding the vertex itself), another one is  $\tilde{C}_B(v) = \frac{C_B(v) - \min C_B}{\max C_B - \min C_B}$

Computation time can be very long,  $O(n_V^3)$ , but one can reach  $O(n_V n_E)$ , using Brandes (2001, [A faster algorithm for betweenness centrality](#))

$v^* = \text{argmax}\{C_B(v)\}$  is seen as the controller of information flow

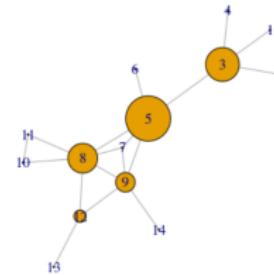
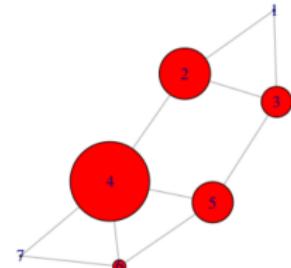
# Networks with R

```
1 > A = get.adjacency(g1, sparse=FALSE)
2 > ng = network::as.network.matrix(A)
3 > sna::gplot.target(ng, betweenness(ng), circ.col="skyblue", vertex.col=c("blue", rep("red", 5), "yellow"), edge.col="darkgray")
4 > betweenness(ng)
5 [1] 0.000 6.667 4.000 10.333
     5.333 1.667 0.000
```

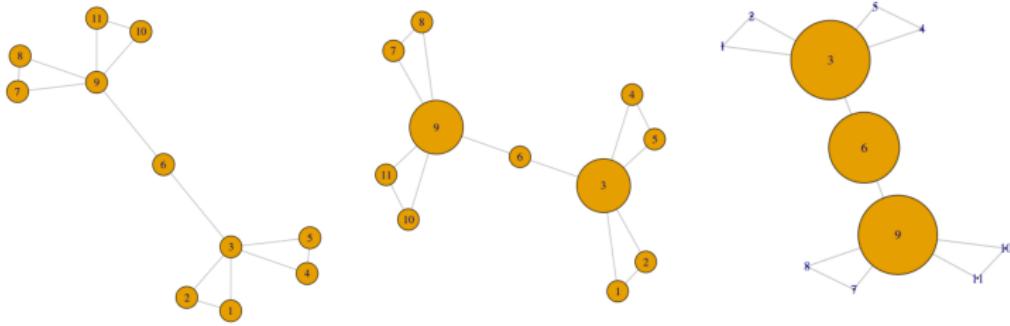


# Networks with R

```
1 > plot(g, vertex.size = betweenness  
       (ng))
```



# Networks & Centrality



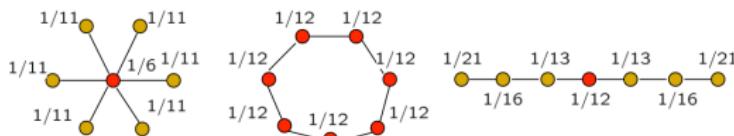
## Networks & Centrality

Closeness is based on the length of the average shortest path between a vertex and all vertices in the graph,

$$C_C(v) = \frac{1}{\sum_y \text{dist}(y, v)}.$$

where  $\text{dist}$  is the geodesic distance, and normalized Closeness Centrality is

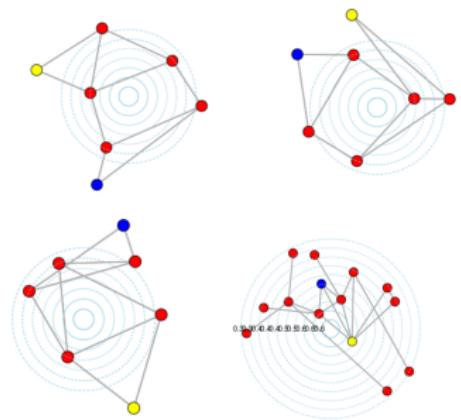
$$C_{C'}(v) = \frac{C_C(v)}{n_V} \text{ or } \frac{C_C(v)}{n_V - 1}$$



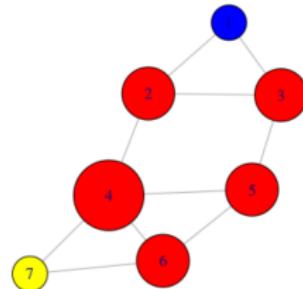
Computation time is  $O(n_V^2 \log n_V + n_V n_E)$  on sparse network.  
 $v^* = \text{argmax}\{C_C(v)\}$  is the **most approachable vertex**.

# Networks with R

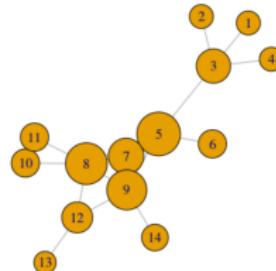
```
1 > A = get.adjacency(g1, sparse=FALSE)
2 > ng = network::as.network.matrix(A)
3 > sna::gplot.target(ng, closeness(ng), circ.col="skyblue", vertex.col=c("blue", rep("red", 5), "yellow"), edge.col="darkgray")
4 > closeness(ng)
5 [1] 0.5000 0.6667 0.6000 0.7500
     0.6667 0.6000 0.5000
```



# Networks with R

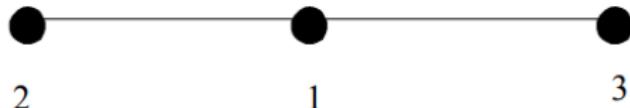


```
1 > plot(g, vertex.size=closeness(ng))
```



## Networks with R

Consider a network with 3 nodes and adjacency matrix



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

```
1 > A=matrix(c(0,1,1,1,0,0,1,0,0),3,3)
2 > A%*%A
3   [,1]  [,2]  [,3]
4 [1,]    2     0     0
5 [2,]    0     1     1
6 [3,]    0     1     1
```

## Networks with R

One can prove that

$$\mathbf{A}^{2k} = \begin{pmatrix} 2^k & 0 & 0 \\ 0 & 2^{k-1} & 2^{k-1} \\ 0 & 2^{k-1} & 2^{k-1} \end{pmatrix} \text{ and } \mathbf{A}^{2k+1} = \begin{pmatrix} 0 & 2^k & 2^k \\ 2^k & 0 & 0 \\ 2^k & 0 & 0 \end{pmatrix}$$

```
1 > eigen(A)
2 $values
3 [1] 1.414214 0.000000 -1.414214
4
5 $vectors
6 [,1] [,2] [,3]
7 [1,] -0.7071068 0.0000000 0.7071068
8 [2,] -0.5000000 -0.7071068 -0.5000000
9 [3,] -0.5000000 0.7071068 -0.5000000
```

## Networks with R

Eigenvalues are  $\{\sqrt{2}, 0, -\sqrt{2}\}$ . Eigenvector with positive components associated to the largest eigenvector (see [Perron-Frobenius theorem](#)) is

$$\mathbf{x} = \frac{1}{2} (\sqrt{2}, 1, 1)$$

```
1 > eigen_centrality(g)
2 $vector
3      1           2           3
4 1.0000000 0.7071068 0.7071068
5
6 $value
7 [1] 1.414214
```

## Networks & Centrality

This eigenvalue centrality score is related to the following (recursive) equation

$$C_E(v) \propto \sum_{u \in M(v)} C_E(u) = \sum_{u \in E} A_{v,u} C_E(u)$$

where  $M(v)$  is the set of neighbors of  $v$ , or equivalently

$$\mathbf{AC}_E \propto \mathbf{C}_E$$

which is an eigenvector equation,  $\mathbf{AC}_E = \lambda \mathbf{C}_E$  for some  $\lambda$ .

## Networks & Centrality

**Theorem [Perron–Frobenius]** Let  $A = (a_{ij})$  be an  $n \times n$  matrix with positive entries,  $a_{ij} > 0$  for  $1 \leq i, j \leq n$ , then

- There is a positive real number  $r$ , called the Perron–Frobenius eigenvalue, such that  $r$  is an eigenvalue of  $A$  and any other eigenvalue (possibly, complex) is strictly smaller than  $r$  in absolute value
- There exists an eigenvector  $\mathbf{x}$  of  $A$  with eigenvalue  $r$  such that all components of  $\mathbf{x}$  are positive
- There are no other positive eigenvectors except positive multiples of  $\mathbf{x}$

**Proof** Meyer (2000, *Matrix analysis and applied linear algebra*)  
Hence,  $\lambda$  is Perron–Frobenius eigenvalue, and  $C_E$  is the (normalized) eigenvector

## Networks & Centrality

Other related measures can be considered.

Katz status index is a weighted count of all paths coming to the node, attenuated with factor  $\beta$ .

Define the vector of Bonacich centralities of parameter  $\beta$ , for  
 $\mathbf{u} \in \mathbb{R}_+^{n_V}$

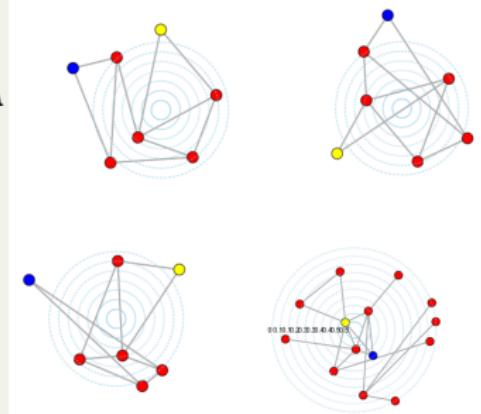
$$b_1(\mathbf{G}, \beta) = \sum_{h \geq 1} \beta^h \mathbf{A}^h \mathbf{1} = ([\mathbb{I} - \beta \mathbf{A}]^{-1} - \mathbb{I}) \mathbf{1}$$

if  $\beta$  is smallest than the inverse of the largest eigenvalue, see see  
Bonacich (1987), or

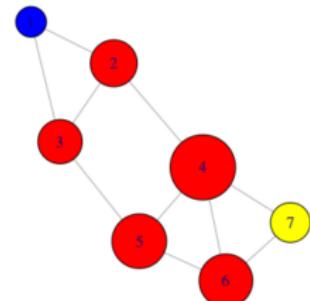
$$b_1(\mathbf{G}, \beta) = \beta [\mathbb{I} - \beta \mathbf{A}]^{-1} \mathbf{A} \mathbf{1}$$

# Networks with R

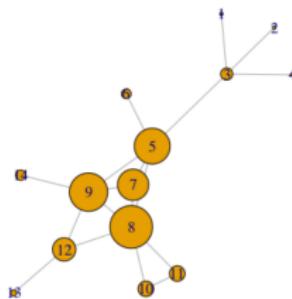
```
1 > A = get.adjacency(g1, sparse=FALSE)
2 > ng = network::as.network.matrix(A)
3 > sna::gplot.target(ng, evcent(ng),
        circ.col="skyblue", vertex.col=
        c("blue", rep("red", 5), "yellow"), edge.col="darkgray")
4 > evcent(ng)
5 [1] 0.2358069 0.3622484 0.3416056
     0.5038536 0.4215946 0.4129477
     0.3071490
```



# Networks with R



```
1 > plot(g, vertex.size=evcent(ng))
```



## Networks & Centrality

Centrality measures we've seen use the shortest path between pairs of vertices.

Katz centrality measures considers the total number of paths between pairs of vertices.

$$C_K(i) = \sum_{k=1}^{\infty} \sum_{j=1}^{n_V} \alpha^k (\mathbf{A}^k)_{ji}$$

with some attenuation coefficient  $\alpha$ .

$\alpha$  has to be smaller than the inverse of the largest eigenvalue of the adjacency matrix  $\mathbf{A}$ .

One can also write  $C_K = ([\mathbb{I} - \alpha \mathbf{A}^T]^{-1} - \mathbb{I}) \mathbf{1}$

Can be related to random walks on a graph : from vertice  $v$ , with probability  $\alpha$  move to a neighbor of  $v$ , and with probability  $1 - \alpha$  move to any vertice in  $V$ .

## Networks & Centrality

Score of a page is proportional to the sum of the scores of pages linked to it

Define  $\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$  where  $\mathbf{D} = \text{diag}(\mathbf{d})$  is the diagonal matrix of degrees (out-degree for internet pages), and consider the following update rule

$$R^{(k)}(v) = \sum_{u \in M(v)} \frac{1}{d_u} R^{(k-1)}(u) \text{ i.e. } \mathbf{R}^{(k)} = \mathbf{P}^T \mathbf{R}^{(k-1)}$$

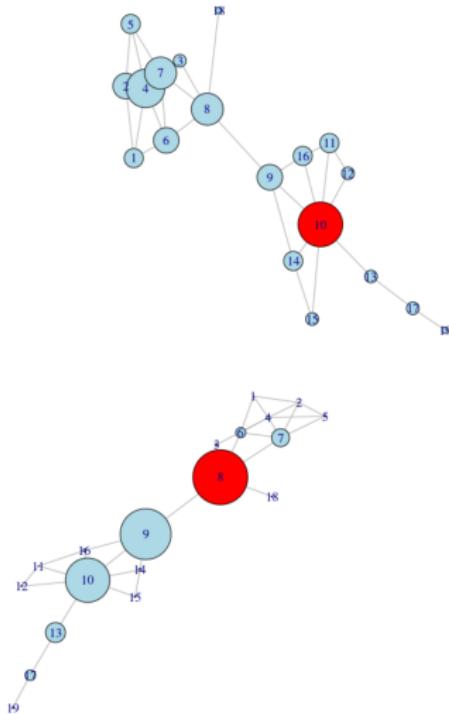
initialized with  $\mathbf{R}^{(0)} = \mathbf{1}/N_V$ . Or consider some scaled version

$$\mathbf{R}^{(k)} = \tilde{\mathbf{P}}^T \mathbf{R}^{(k-1)} \text{ where } \alpha \tilde{\mathbf{P}} + (1 - \alpha) \frac{\mathbf{1}\mathbf{1}^T}{N_V}$$

where the scaling factor  $\beta$  is typically around 0.85  
power iteration converges to the dominant eigenvector of  $\tilde{\mathbf{P}}$ .

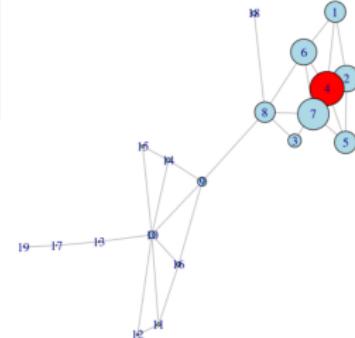
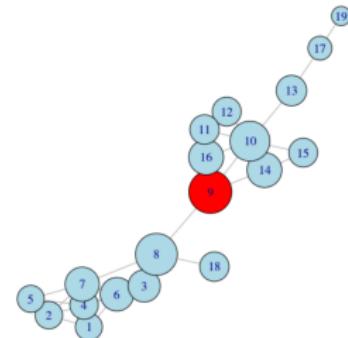
# Networks with R

```
1 > A <- get.adjacency(g, sparse=FALSE)
2 > library(network)
3 > ng <- network::as.network.matrix(A)
4 > library(sna)
5 > which.max(sna::degree(ng))
6 [1] 13
7 > which.max(sna::betweenness(ng))
8 [1] 8
```



# Networks with R

```
1 > A <- get.adjacency(g, sparse=FALSE)
2 > library(network)
3 > ng <- network::as.network.matrix(A)
4 > library(sna)
5 > which.max(sna::closeness(ng))
6 [1] 9
7 > which.max(sna::evcent(ng))
8 [1] 3
```



None is wrong, they just focus on different notions of 'importance'

## Networks & Centrality

Need a metric on the space of networks. Let  $G_1 = (E_1, V)$  and  $G_2 = (E_2, V)$

$$\Delta(G_1, G_2) = \sum_{e \in E_1 \cup E_2} |\mathbf{1}_{e \in E_1} - \mathbf{1}_{e \in E_2}|$$

A centrality measure  $C$  is said to be **stable** if for any  $v \in V$  and any graphs  $G_1$  and  $G_2$ ,

$$|C^{G_1}(v) - C^{G_2}(v)| \leq \kappa \Delta(G_1, G_2)$$

for some constant  $\kappa$  (see Lipschitz continuity).

Idea: the importance of nodes should be robust to small perturbations in the graph.

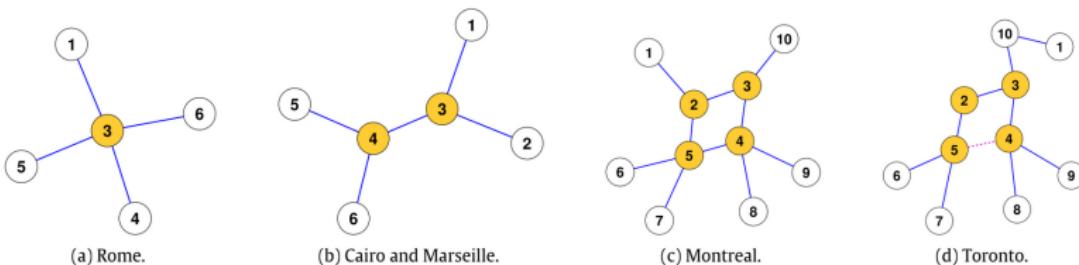
For instance, one can prove that for the degree

$$|C_D^{G_1}(v) - C_D^{G_2}(v)| \leq \Delta(G_1, G_2)$$

thus, **degree centrality is a stable measure.**

# Centrality, Vulnerability & Security

See Worldwide Metro Networks, from Wang *et al.* (2017,  
Multi-criteria robustness analysis of metro networks)



Defined several measures of robustness of metro networks around the world

# Centrality, Vulnerability & Security

See Worldwide Metro Networks, from Wang et al. (2017,  
Multi-criteria robustness analysis of metro networks)

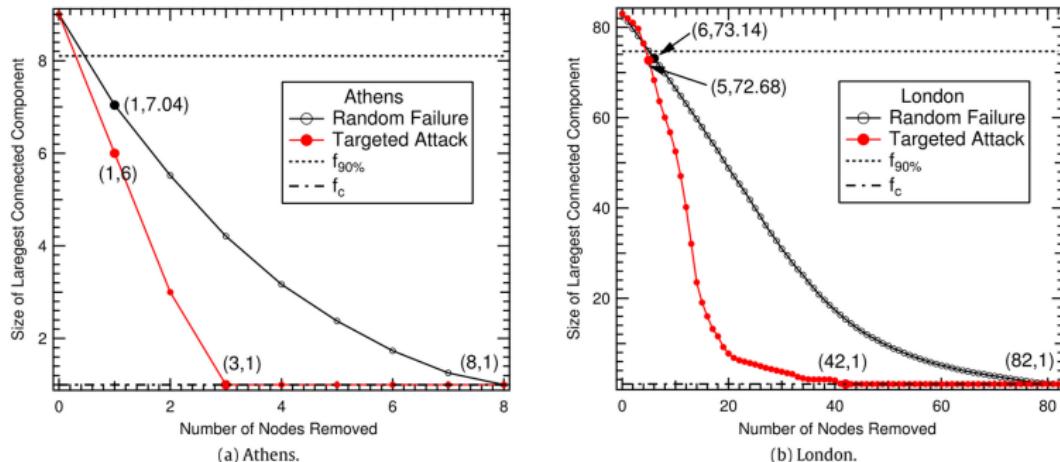


Fig. 4. Critical thresholds in metro networks under nodes removal.

## Inference & Networks

Recall **classical statistical tools**: consider a sample of observations  $\mathbf{x} = \{x_1, \dots, x_n\}$  from i.i.d. random variables  $X_1, \dots, X_n$ , with distribution  $F$ .

Assume that  $F \in \mathcal{F} = \{F_\theta, \theta \in \Theta\}$

Consider some point estimator  $\hat{\theta} = s(x_1, \dots, x_n)$

Estimator can be seen as a random variable  $\hat{\theta} = s(X_1, \dots, X_n)$

The bias is  $\mathbb{E}[\hat{\theta}] - \theta$  and its standard-error  $\text{Var}[\hat{\theta}]^{1/2}$ . The mean squared error is

$$\text{mse}[\hat{\theta}] = \mathbb{E}([\hat{\theta} - \theta]^2) = \text{bias}[\hat{\theta}]^2 + \text{Var}[\hat{\theta}]$$

## Inference & Networks

Survey sampling is a standard tool in socio-economic studies

Network sampling designs provide an alternative

see McCormick, He, Kolaczyk & Zheng (2012, Surveying  
hard-to-reach groups through sampled respondents in a social  
network)

Let  $G = (V, E)$  denote a network, and let  $\theta(G)$  denote a statistics  
of interest, e.g.

- the number of nodes  $n_V$
- the number of connexions  $n_E$
- the degree  $d_v$  of a node  $v$

using a sampled subgraph of  $G$  is maybe not a great idea...

## Inference & Networks

Measurements are usually on a portion of the complex population. Here, we cannot observe  $G = (V, E)$  but only  $G^* = (V^*, E^*)$ , a subgraph of  $G$ .

Consider some statistics of interest on the graph,  $\theta(G)$ .

Can we consider the natural plug-in estimator  $\hat{\theta} = \theta(G^*)$  ?

Consider the case where  $\theta(G)$  is the average degree,

$$\theta(G) = \frac{1}{n_V} \sum_{v \in V} d_v. \text{ How could we sample ?}$$

- sample  $n$  vertices  $V_n^* = v_1, \dots, v_n$  (without replacement)
  - for each  $v_k \in V_n^*$  observe incident edges  $(v_k, v) \in E$
  - observe edges only when vertices are in  $V^*$ ,  $v_{k_1} \in V_n^*$  and  $v_{k_2} \in V_n^*$ .

$$\text{Set } \hat{\theta} = \frac{1}{n} \sum_{v \in V_n^*} d_v^*$$

## Inference & Networks

The first technique is fine,  $\hat{\theta} \sim \theta(G)$  but the second one under-estimate  $\theta(G)$  (actually,  $\hat{\theta} \sim n\theta(G)/n_V$ ).

Necessary to incorporate effects of random sampling.

Consider there the mean  $\mu$  of  $y_i$ 's.

Let  $\pi_i$  denote the probability to have individual  $i$ , in the sample.

$$\mathbb{E}[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^m \pi_i y_i \neq \mu \text{ if } \pi_i \neq \frac{n}{m}$$

See Horvitz-Thompson's estimator (see Horvitz & Thompson (1952)): unequal probability sampling necessitates unequal weights when averaging : use inclusion probabilities as weights.

## Inference & Networks

An unbiased estimator of  $\mu$  is  $\tilde{\mu}_{HT} = \frac{1}{m} \sum_{i \in \text{sample}} \frac{y_i}{\pi_i}$ ,

$$\mathbb{E}[\tilde{\mu}_{HT}] = \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^n \frac{y_i}{\pi_i} \underbrace{\mathbf{1}(i \text{ is in sample})}_{=Z_i} \right] = \frac{1}{m} \sum_{i=1}^n \frac{y_i}{\pi_i} \underbrace{\mathbb{E}[Z_i]}_{=\pi_i} = \mu.$$

Further

$$\text{Var}[\tilde{\mu}_{HT}] = \frac{1}{n^2} \sum_{i,j=1}^n y_i y_j \left( \frac{\pi_{i,j}}{\pi_i \pi_j} - 1 \right)$$

The main problem here is to compute  $\pi$ 's. Furthermore

$$\text{Var}[\tilde{\mu}] = \frac{1}{m^2} \sum_{i=1}^m \left( \frac{1 - \pi_i}{\pi_i^2} \right) y_i^2 + \frac{1}{m^2} \sum_{i \neq j=1}^m \left( \frac{\pi_{i,j} - \pi_i \pi_j}{\pi_i \pi_j} \right) y_i y_j$$

where  $\pi_{i,j}$  is the probability to have  $i$  and  $j$  in the sample.

## Inference & Networks

with simple random sampling (without replacement)

$$\pi_{i,j} = \frac{n(n-1)}{m(m-1)}.$$

but need to know  $m$ .

Use the [mark and recapture estimator](#), see Chao *et al.* (2001, [The applications of capture-recapture models to epidemiological data](#)) to estimate  $m$ . Consider two samples  $\mathcal{S}_1$  and  $\mathcal{S}_2$ .

Step 1: mark all units in sample  $\mathcal{S}_1$  of size  $n_1$

Step 2: create a sample  $\mathcal{S}_2$  of size  $n_2$

$$\hat{m} = \frac{\text{Card}(\mathcal{S}_1) \cdot \text{Card}(\mathcal{S}_2)}{\text{Card}(\mathcal{S}_1 \cap \mathcal{S}_2)}$$

Complicated to estimate if the population is  $V$  (vertices) or  $E$  (edges).

## Inference & Networks

Idea: take  $n$  vertices, set  $V_n$ , and observe all edges in the subgraph induced by  $V_n$ ,

$$\pi_i = \frac{n}{n_V} \text{ and } \pi_{i,j} = \frac{n(n-1)}{n_V(n_V-1)}$$

for vertices inclusion. See friendship paradox

## Inference & Networks

Idea: take  $n$  edges, set  $E_n$ , and observe all vertices incident to edges in  $E_n$ ,

$$\pi_i = \frac{n}{n_V} \text{ and } \pi_{i,j} = \frac{n(n-1)}{n_V(n_V-1)}$$

See phone calls

Vertices inclusion probabilities are here

$$\pi_i = \mathbb{P}[\text{vertex } i \text{ is sampled}] = 1 - \frac{\binom{n_E - d_i}{n}}{n_E} n$$

if  $n \leq n_E - d_i$  (1 otherwise).

# Inference & Networks

## Inference & Networks

also called **Snowball Sampling**

Idea: start with an initial vertex sample  $V_0$  and observe all incident edges, as well as all vertices sharing those edges. Iterate to have  $n$  edges.

## Inference & Networks

### Link Tracing

Idea: start with an initial vertex sample  $V_0$  and observe *some* incident edges (called 'links'), as well as all vertices sharing those edges. Iterate to have  $n$  edges.

### Path Sampling

Select a set of source edges, and target sources, and consider some shortest paths from each source to all targets. Called **Traceroute sampling**. One can prove that

$$\pi_i \sim 1 - \left(1 - \frac{n_s}{n_E} \frac{n_t}{n_E}\right) \exp\left(-\frac{n_s}{n_E} \frac{n_t}{n_E} b_i\right)$$

and

$$\pi_{i,j} \sim 1 - \exp\left(-\frac{n_s}{n_E} \frac{n_t}{n_E} b_{i,j}\right)$$

where  $b_i$  is the betweenness centrality of vertice  $i$  and  $b_{i,j}$  is the betweenness centrality of edge  $(i,j)$ .

## Inference & Networks

see Dall'Asta *et al.* (2006, [Vulnerability of weighted networks](#)) for more details.

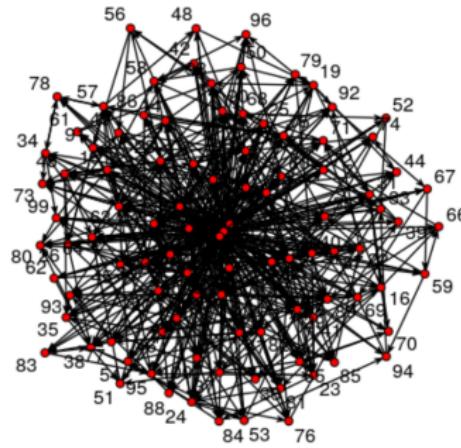
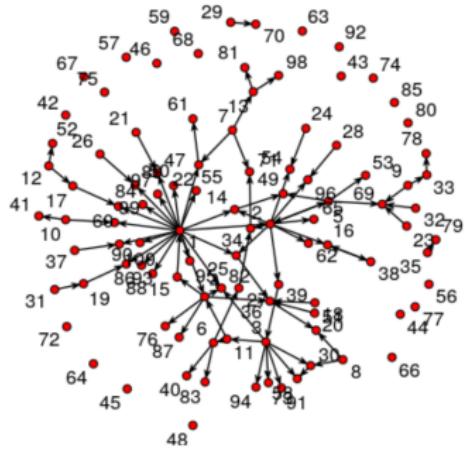
## Inference & Networks

emotion/behaviors of individuals being influenced by the group(s) to which they belong

see Christakis & Fowler (2009, [Connected : The Surprising Power of Our Social Networks and How They Shape Our Lives](#) )

challenge of distinguishing social contagion from ordinary homophily

## Inference & Networks



## Networks with R

Do friends of peoples tend to be friends together?

The cohesion can be related to density, clustering, connectivity or transitivity (*the friend of your friend is likely to be your friend*) etc.

The density of a sub-graph  $H \subset G$  is defined as

$$\text{density}(H) = \frac{2n_{E_H}}{n_{V_H}(n_{V_H} - 1)}$$

where  $H = (V_H, E_H)$ , which is the number of edges over the maximum number of edges (based on the number of vertices).

```
1 > graph.density(g1)
2 [1] 0.4761905
3 > graph.density(g2)
4 [1] 0.1978022
```

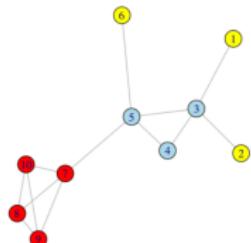
Usually, we consider for  $H$  the neighborhood of a vertice  $v$ .

## Clusters & Networks

A **clique** is a complete subgraph of  $G$ .

Note that large cliques are rare, since a single missing edge destroys the property.

A sufficient condition to have a clique of size  $n$  is



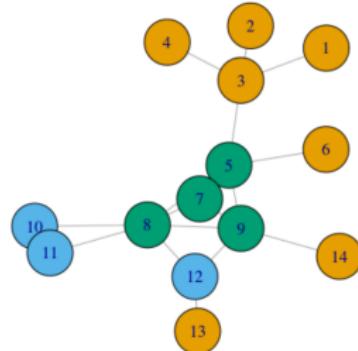
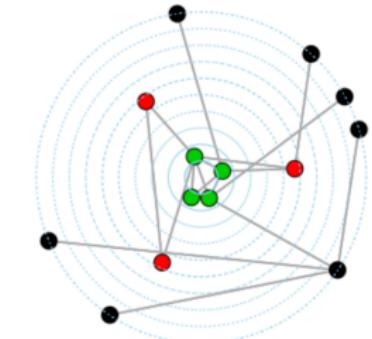
$$n_E > \frac{n_V^2}{2} \frac{n-2}{n-1}$$

A  **$k$ -core**  $G^*$  is a subgraph of  $G$  such that  $d_{v^*} \geq k$  for all  $v^* \in V^*$ , and  $G^*$  is maximal.

Here, degrees are computed on the subgraph  $G^*$ .

# Networks with R

```
1 A = get.adjacency(g2, sparse=FALSE)
2 ty= network::as.network.matrix(A)
3 cores = graph.coreness(g2)
4 sna::gplot.target(ty, cores, circ.
  col="skyblue", usearrows =
  FALSE, vertex.col=cores, edge.
  col="darkgray")
5 V(g2)$color = cores
6 plot(g2)
```



## Clusters & Networks

Note that **clustering** is related (heuristically) on cycles among 3 vertices (triangles).

$\tau_\Delta(v)$  is the number of triangles in  $G$  that contain  $v$

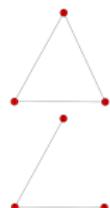
$\tau_\Lambda(v)$  is the number of connected triplets in  $G$  with two edges incident to  $v$

One can define an aggregate local clustering coefficient,

$$\text{cluster}(v) = \frac{\tau_\Delta(v)}{\tau_\Lambda(v)}$$

Let  $V_2$  denote the subset  $\{v \in V, d_v \geq 2\}$ , then set

$$\text{cluster}(G) = \frac{1}{n_{V_2}} \sum_{v \in V_2} \text{cluster}(v)$$



the clustering coefficient, and a transitivity coefficient

$$\text{transitivity}(G) = \frac{\sum_{v \in V_2} \text{cluster}(v) \cdot \tau_\Lambda(v)}{\sum_{v \in V_2} \tau_\Lambda(v)}$$

# Networks with R

```
1 > transitivity(g1)
2 [1] 0.45
3 > transitivity(g2)
4 [1] 0.3673469
```

One can also look for [network partitioning](#), i.e. community detection

- vertices in a group should be well connected among themselves
- vertices in different groups should be well separated

Hierarchical clustering and spectral partitioning

## Clusters & Networks

We seek a segmentation in natural subsets, i.e. a partition of the set of vertices.

A subset of vertices is said to be **cohesive** if

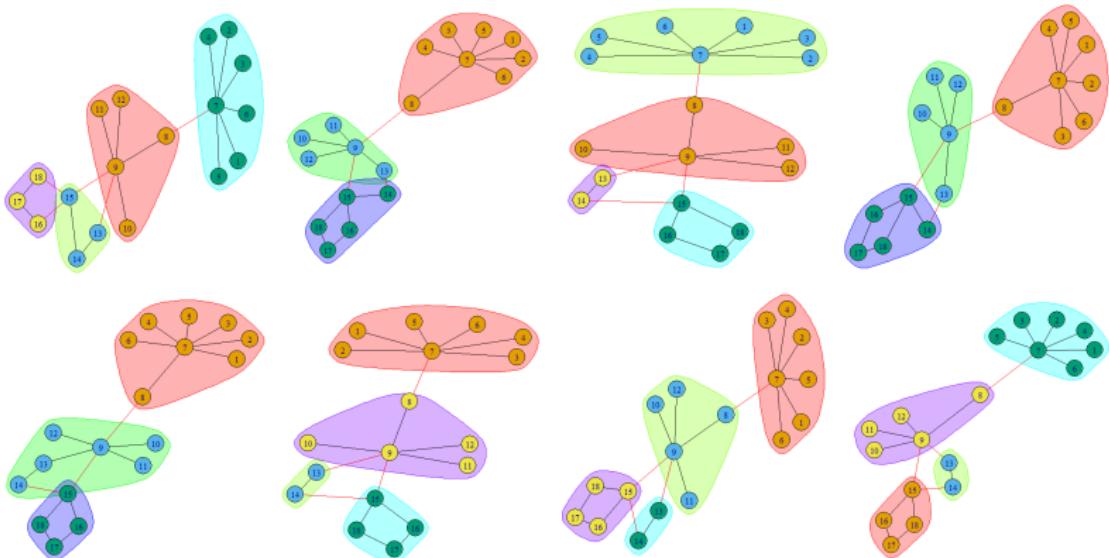
- vertices are well connected among themselves (within a subset)
- vertices are well separated from the remaining vertices (between subsets)

i.e. given a partition  $\{C_1, \dots, C_k\}$

- **between**:  $E(C_k, C_{k'})$  set of connecting vertices in  $C_k$  to vertices in  $C_{k'}$  is relatively small in size compared to
- **within**  $E(C_k)$  and  $E(C_{k'})$ , edges connecting vertices within subsets

**Application** community detection in social media

# Clusters & Networks



Hierarchical clustering (greedy approach, since we iteratively modify partitions)

Can be either forward or backward,

- successive coarsening of partition, by merging
- successive refinement of partition, by splitting

with two extremal partitions,  $V$  and  $\{\{v_1\}, \dots, \{v_k\}\}$ .

Dendrogram-based representation.

(1) Need a cost measure to quantify (dis)similarity between pairs of vertices.

(2) Need to aggregate (dis)similarities on subsets  $C_k$  and  $C_{k'}$

- classical single linkage
- common linkage

## Clusters & Networks

- Ward's method

Euclidean Similarity

$$ds_{i,j} = \sqrt{\sum_{k \neq i,j} (A_{i,k} - A_{j,k})^2}$$

Neighborhood-based Similarity

$$ds_{i,j} = \frac{\Delta_{i,j}}{d_{N_k} + d_{N_k}}$$

where  $\Delta_{i,j}$  is the number of neighbors shared by  $i$  and  $j$

## Clusters & Networks

One can use modularity, from Newman (2004, Modularity and community structure in networks)

Heuristically, we choose a partition  $C$  that deviates most in modularity from what is expected at random from a collection of graphs with the same degree distribution as  $G$ .

Consider a partition  $\mathcal{C} = \{C_1, \dots, C_k\}$ . Given two clusters  $i$  and  $j$ , define  $f_{i,j}(\mathcal{C})$  the fraction of edges connecting vertices in  $C_i$  to vertices in  $C_j$ .

Modularity of  $\mathcal{C}$  is

$$\text{mod}(\mathcal{C}) = \sum_k [f_{k,k}(\mathcal{C}) - f_{k,k}^*]$$

where  $f_{k,k}^*$  is the expected value of  $f_{k,k}$  under some random edge assignment.

## Clusters & Networks

One can also consider **spectral partitioning** based on spectral (eigenvalues) of a matrix related to network  $G$ , e.g. adjacency matrix  $\mathbf{A}$  or Laplacian matrix  $\mathbf{L}$ .

On  $\mathbf{A}$ , calculate  $n_V$  eigenvalues and eigenvectors,

$$(\lambda_i, \mathbf{u}_i), \text{ with } \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_{n_V}.$$

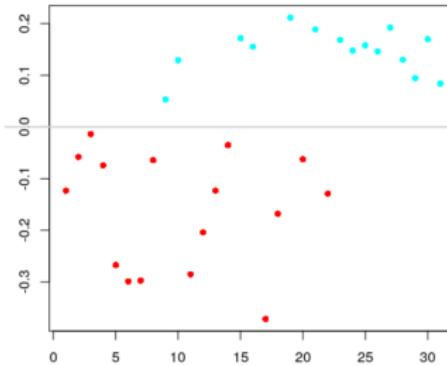
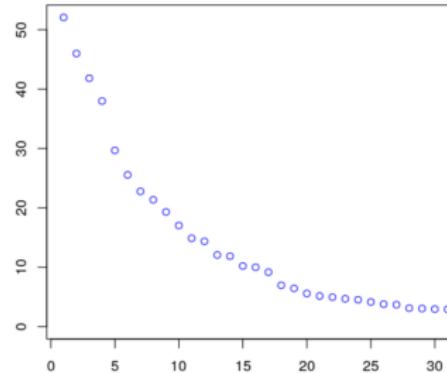
Start with the largest (absolute value) eigenvalues, vertices corresponding to particularly large positive or negative entries, in conjunction with their immediate neighbors, are declared to be a cluster

On  $\mathbf{L} = \mathbf{D} - \mathbf{A}$ , recall that  $\mathbf{x}^T \mathbf{L} \mathbf{x} = \sum_{i,j} (x_i - x_j)^2$

The closer  $\mathbf{x}^T \mathbf{L} \mathbf{x}$  is to zero, the more similar the elements of  $\mathbf{x}$  at adjacent vertices in  $V$ . So the Laplacian provides some sense of the 'smoothness' of functions  $\mathbf{x}$  on the network.

# Networks with R

```
1 > k.lap <- graph.laplacian(karate)
2 > eigk <- eigen(k.lap)
3 > plot(eigk$values, col="blue")
4 > f.vec <- eigk$vectors[, 33]
5 >
6 > faction <- get.vertex.attribute(
    karate, "Faction")
7 > f.colors <- as.character(length(
    faction))
8 > f.colors[faction == 1] <- "red"
9 > f.colors[faction == 2] <- "cyan"
10 > plot(f.vec, pch=16, col=f.colors)
11 > abline(h=0, 0)
```



## Walktrap

This algorithm finds densely connected subgraphs by performing random walks. The idea is that random walks will tend to stay inside communities instead of jumping to other communities.

see [Pons & Latapy \(2005\)](#)

## Edge Betweenness

This algorithm is the [Newman & Girvan \(2003\)](#) algorithm. It is a divisive algorithm where at each step the edge with the highest betweenness is removed from the graph. For each division you can compute the modularity of the graph. At the end, choose to cut the dendrogram where the process gives you the highest value of modularity.

## Clusters & Networks

### Fastgreedy

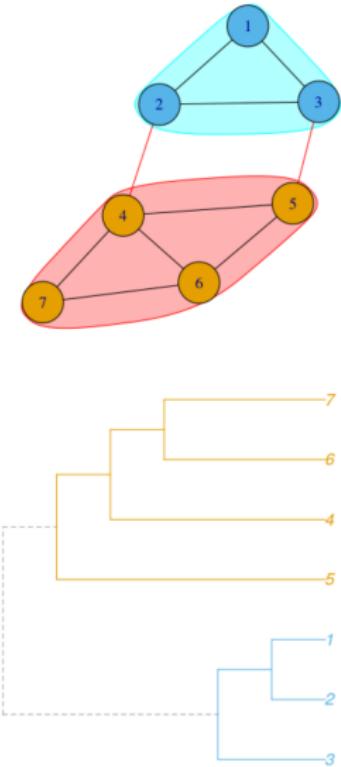
This algorithm is the [Clauset, Newman & Moore \(2004\)](#) algorithm. In this case the algorithm is agglomerative. At each step two groups merge. The merging is decided by optimising modularity. This is a fast algorithm, but has the disadvantage of being a greedy algorithm. Thus, it might not produce the best overall community partitioning, although I find it useful and accurate.

### Spin-glass

This algorithm uses a spin-glass model ([Sherrington & Kirkpatrick \(1975\)](#)) and simulated annealing to find the communities inside a network. See also [Reichardt & Bornholdt \(2006\)](#)

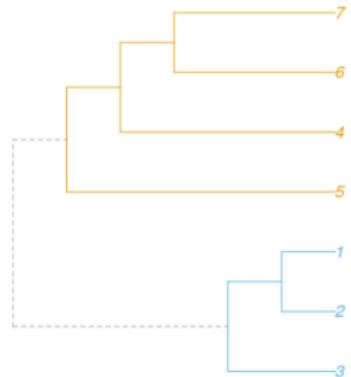
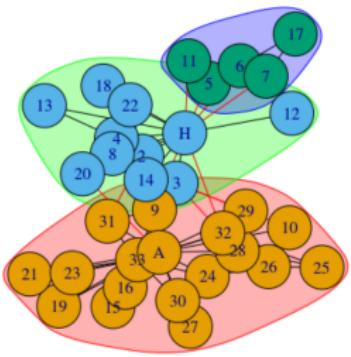
# Networks with R

```
1 > kc = fastgreedy.community(g1)
2 > length(kc)
3 [1] 2
4 > sizes(kc)
5 Community sizes
6 1 2
7 4 3
8 > membership(kc)
9 1 2 3 4 5 6 7
10 2 2 2 1 1 1 1
11 > plot(kc,g1)
12 > dendPlot(kc,mode="phylo")
```



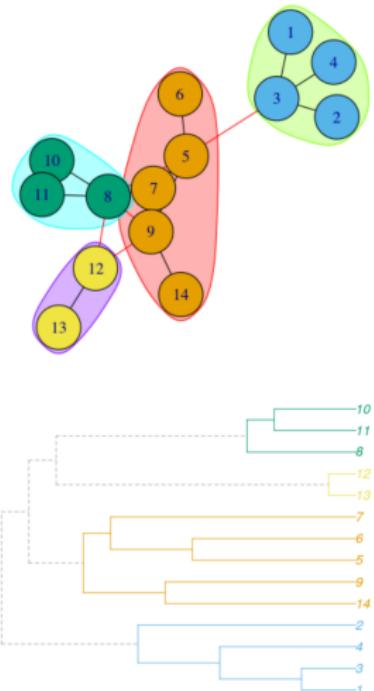
# Networks with R

```
1 > kc <- fastgreedy.community(karate)
2 > length(kc)
3 [1] 3
4 > sizes(kc)
5 Community sizes
6 1 2 3
7 18 11 5
8 > plot(kc,karate)
9 > dendPlot(kc,mode="phylo")
```



# Networks with R

```
1 > kc <- fastgreedy.community(g2)
2 > length(kc)
3 [1] 4
4 > sizes(kc)
5 Community sizes
6 1 2 3 4
7 5 4 3 2
8 > membership(kc)
9  1   3   2   4   5   6   7   8   9   10  11  12
10    13  14
11  2   2   2   2   1   1   1   3   1   3   3   4
12    4   1
11 > plot(kc,g2)
12 > dendPlot(kc,mode="phylo")
```



## Networks & Small World

*"I read somewhere that everybody on this planet is separated by only six other people. Six degrees of separation. Between us and everybody else on this planet. The president of the United States. A gondolier in Venice. Fill in the names. . . . How every person is a new door, opening up into other worlds. Six degrees of separation between me and everyone else on this planet. But to find the right six people..."* Guare (1990, Six Degrees of Separation)

The **small world** property is mathematically intuitive: if the number of vertices within a given distance of a specific node growth exponentially with the distance, then the average path length increases a  $\log n_V$ . Formally, we talk here about the **average path length**

$$\bar{\ell} = \binom{n_V}{2}^{-1} \sum_{u \neq v} d(u, v) = O(\log n_V)$$

if  $d_v = d$ , if we reach everyone after  $k$  hops  $n_V \sim d^k$  i.e.  
 $k = O(\log n_V)$

## Networks & Small World

*“Each person in the world (at least among the 1.59 billion people active on Facebook) is connected to every other person by an average of three and a half other people. The average distance we observe is 4.57, corresponding to 3.57 intermediaries or ‘degrees of separation’”, see Bharat et. (1990, Three and a half degrees of separation)*

An alternative is based on the harmonic mean, related to graph efficiency

$$\ell^{-1} = \frac{2}{n_V(n_V - 1)} \sum_{i>j} d_{i,j}^{-1}$$

## Networks & Homophily

As seen in the introduction, people have a strong tendency to associate with equals, i.e. homophily or assortative mixing.

Consider some categories (partition)  $\{C_1, \dots, C_k\}$ .

Let  $f_{ij}$  denote the fraction of edges joining vertices of categories  $C_i$  and  $C_j$ . Define

$$f_{i \cdot} = \sum_j f_{ij} \text{ and } f_{\cdot j} = \sum_i f_{ij}$$

The assortativity coefficient is

$$R_A = \frac{\sum_i f_{ii} - \sum_i f_{i \cdot} f_{\cdot i}}{1 - \sum_i f_{i \cdot} f_{\cdot i}}$$

see Newman (2003, Mixing patterns in networks)

Here  $f_{i \cdot} f_{\cdot i}$  the expected fraction of edges joining nodes in  $C_i$ .

## Networks & Homophily

Perfect homophily means  $R_A = R_A^+ = 1$ , the upper bound while the lower bound is

$$R_A^- = -\frac{\sum_i f_i \cdot f_{\cdot i}}{1 - \sum_i f_i \cdot f_{\cdot i}} > -1$$

An alternative is based on

$$R = \frac{\text{trace}(\mathbf{f}) - \|\mathbf{f}\|^2}{1 - \|\mathbf{f}\|^2}$$

# Networks & Homophily

Students sitting in dining halls (men/women, and ethnic groups)

Table Demographics in Adams House Dining Hall  
for two 2 dinners.

Time recorded was 6:30 PM for each night

DIAGRAM 1.1 Dinner May 8, 1997			
BY RACE — A = Asian, B = Black, H = Hispanic, W = White			
BB	WWWW	HBBB	
BB	WWWW	HBBB	
W—	WWWW	W—	
W—	WWWW	W—	
WWW	WWWW	WW—	
WWW	WWWW	W—	
	WWW—	A—	
	WWWW	W—	
	AWBB	WW	
	AWBB	WW	
	WAWA	BB	
	BBWA	BB	
	BBBB	WW	
	BBBB	W—	
WWW	WW	HWW	
WWW	WW	HHH	
B—	WW	A—	AA
W—	W—	W—	AA

Avg. Segregation Level = .89

DIAGRAM 1.2  
Dinner May 8, 1997  
BY SEX — M = Male, F = Female

FF	MMPM	FMMF
FM	MM-M	FFFF
	MMFF	W—
	MMFF	W—
FFF	FFFF	MM—
FFF	FFFF	F—
	FFM—	M—
	FFMM	M—
	FMFM	MF
	MMFF	MF
	FFMM	MM
	MMFF	MM
	MMMM	MM
	MMMM	M—
FFF	FFF	
FFF	FFF	
F—	FF	F—
M—	F—	F—
	FF	FF

Avg. Segregation Level = .81

Population divided in two types, each individuals wants at least a proportion  $t \in [0, 1]$  of his neighbors with his own type (say 50%).

## Networks & Homophily

Unsatisfied individuals move.

There is an externality from the move, with a cascading effect.

Need to define a **neighbourhood**, a **preference** function, a **moving rule** (consider either exchanges or empty space).

Schelling model generates **segregation** beyond what people really want, see Chapter 4 in [Easley & Kleinberg \(2010\)](#).

Fixed and unchanging characteristics (gender, ethnicity) can be highly correlated with mutable characteristics (location).

## Networks & Homophily

Consider a network  $(V, E)$ , with  $n$  edges. Let  $x_e$  be an indicator variable if  $e$  is a cross edge (in blue on the graph).

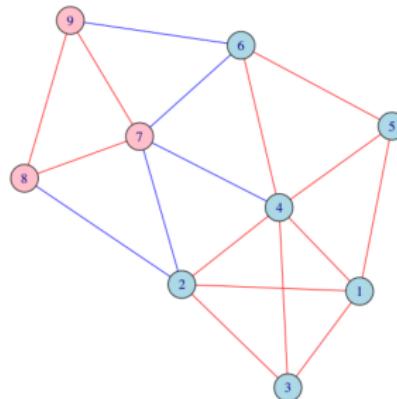
$$\bar{x} = \frac{1}{n_E} \sum_{e \in E} x_e = \frac{5}{18}$$

Let  $p$  denote the proportion of one type (say pink). Here  $p = 1/3$ .

In a random graph,  $X_e$ 's are i.i.d. random variables  $\mathcal{B}(2p(1 - p))$ .

Use a standard  $t$ -test:  $H_0 : \bar{X} > 2p(1 - p)$  against

$H_0 : \bar{X} < 2p(1 - p)$



# Networks with R

```
1 > p=2/3
2 > ne=length(E(g))
3 > nce=5
4 > (x=c(rep(1,nce),rep(0,ne-nce)))
5 [1] 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
6 > t.test(x, alternative = "less", mu=2*p*(1-p))

7
8   One Sample t-test
9

10 data: x
11 t = -1.5342, df = 17, p-value = 0.07169
12 alternative hypothesis: true mean is less than
13   0.4444444
14 95 percent confidence interval:
15     -Inf 0.4667556
16 sample estimates:
17 mean of x
18 0.2777778
```