



# UNIVERSITÀ DEGLI STUDI DI PADOVA

## **Network Science**

A.Y. 23/24

ICT for Internet & multimedia, Data science, Physics of data

# Centrality

Importance of nodes in a network



## Centrality

---

From Wikipedia, the free encyclopedia

*For the statistical concept, see [Central tendency](#).*

In [graph theory](#) and [network analysis](#), indicators of **centrality** identify the most important [vertices](#) within a graph.

Applications include identifying the most influential person(s) in a [social network](#), key infrastructure nodes in the [Internet](#) or [urban networks](#), and [super-spreaders](#) of disease. Centrality concepts were first developed in [social network analysis](#), and many of the terms used to measure centrality reflect their [sociological](#) origin.<sup>[1]</sup> They should not be confused with [node influence metrics](#), which seek to quantify the influence of every node in the network.



[Degree centrality](#) [\[ edit \]](#)

*Main article: [Degree \(graph theory\)](#)*

[PageRank centrality](#) [\[ edit \]](#)

*Main article: [PageRank](#)*

[Betweenness centrality](#) [\[ edit \]](#)

*Main article: [Betweenness centrality](#)*

[Eigenvector centrality](#) [\[ edit \]](#)

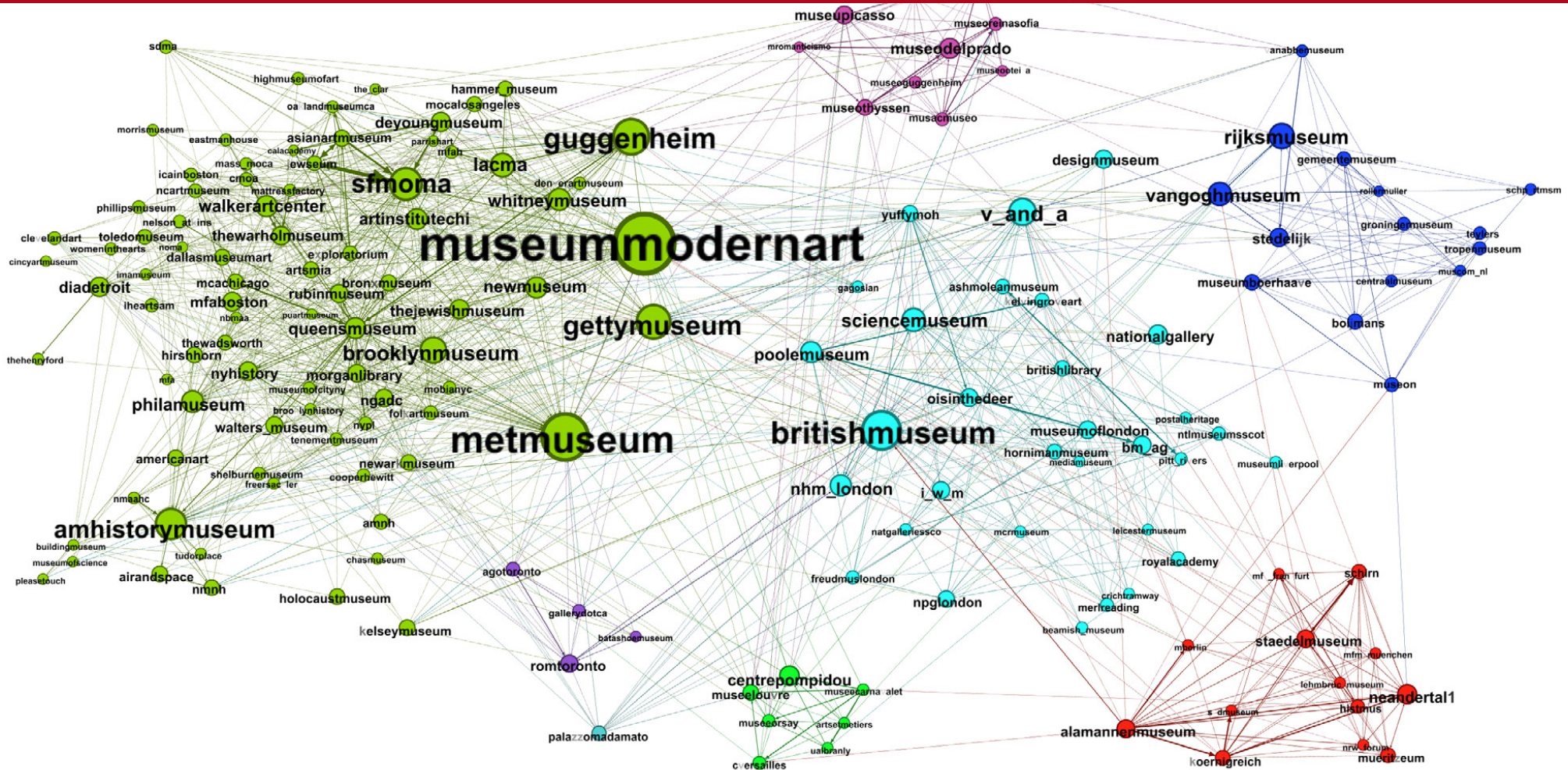
*Main article: [Eigenvector centrality](#)*

[Closeness centrality](#) [\[ edit \]](#)

*Main article: [Closeness centrality](#)*



# An example of node centrality museums network



Can we do this **efficiently**, i.e., by using automatic, reliable, and fast methods?

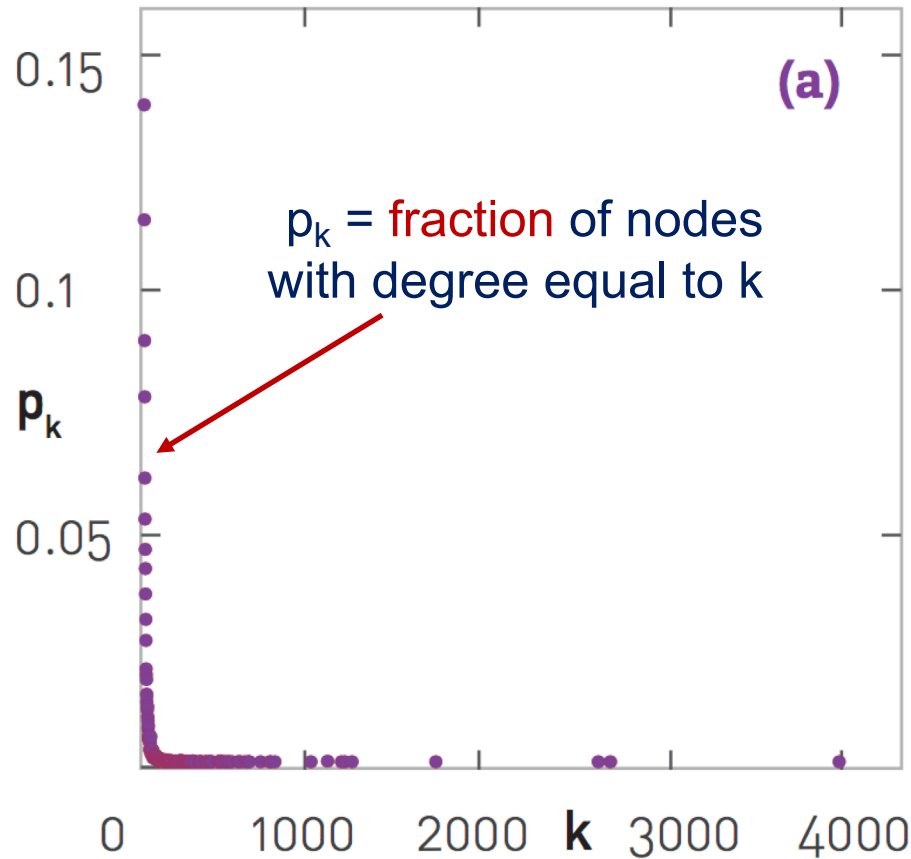
# Degree centrality

Counting the in/out degrees of nodes



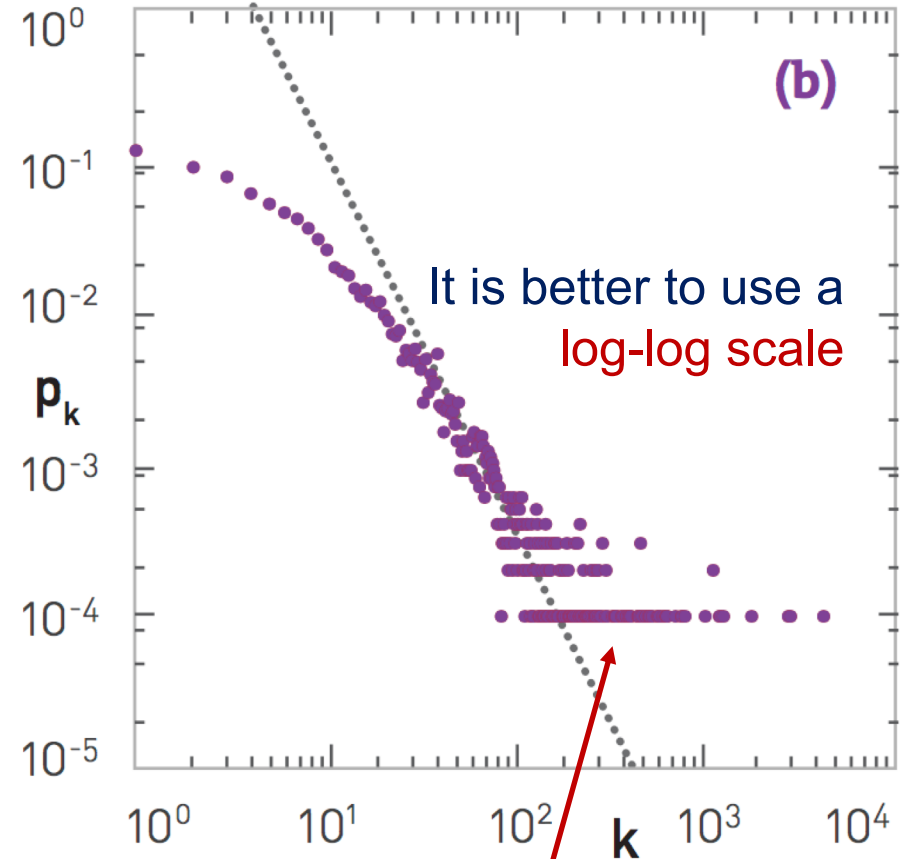
# The degree distribution for an undirected network

~~LINEAR SCALE~~



Wide range for the degree  $k$  !  
Wide range for the probability  $p_k$  !

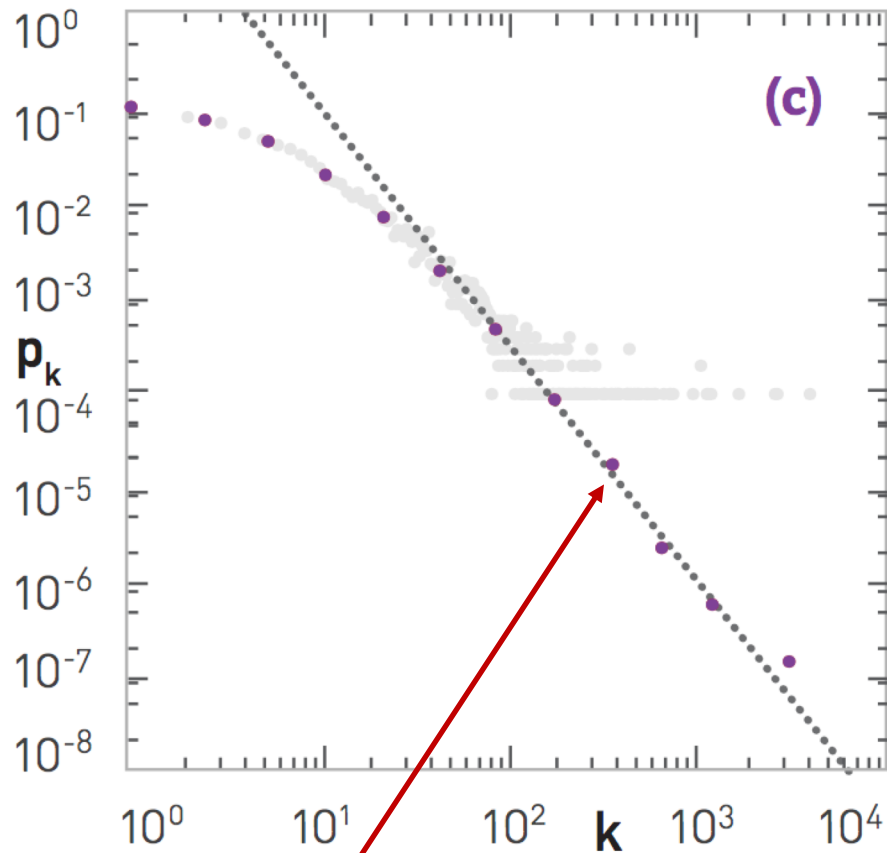
LINEAR BINNING





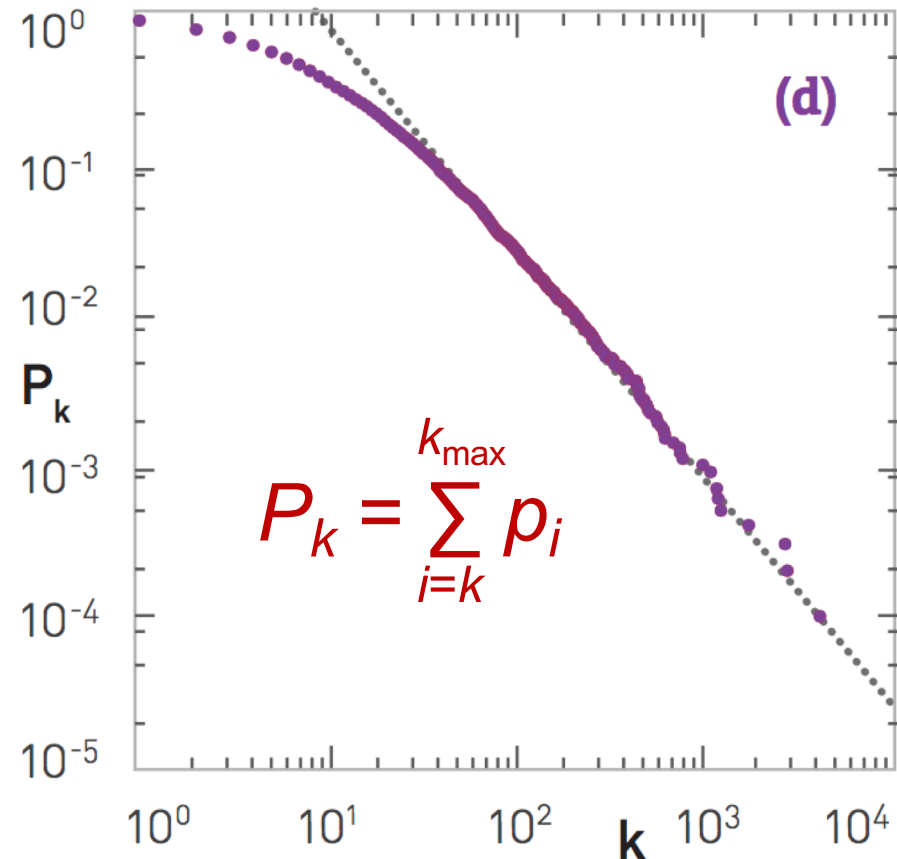
# Alternative log representations for an undirected network

LOG-BINNING



$p_{k_i}$  = fraction of nodes with degree in the range  $[k_i, k_{i+1})$  where  $k_i$  are uniformly distributed in the log-domain,  $k_{i+1} = k_i \cdot \Delta$

CUMULATIVE

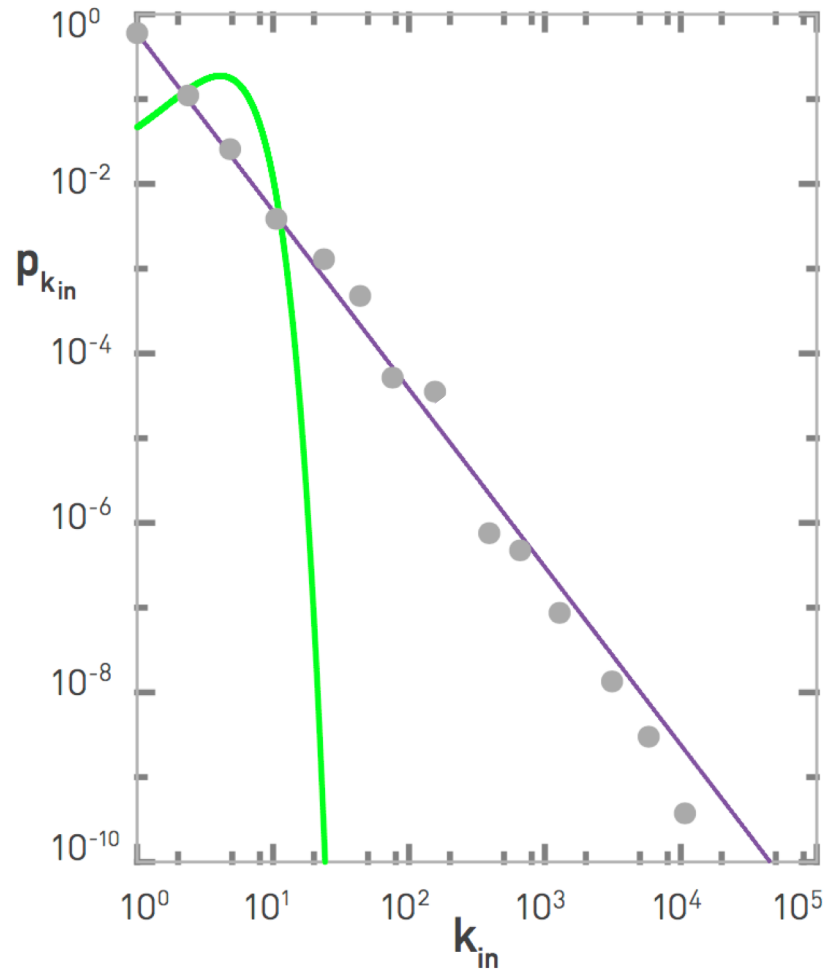


$$P_k = \sum_{i=k}^{k_{\max}} p_i$$

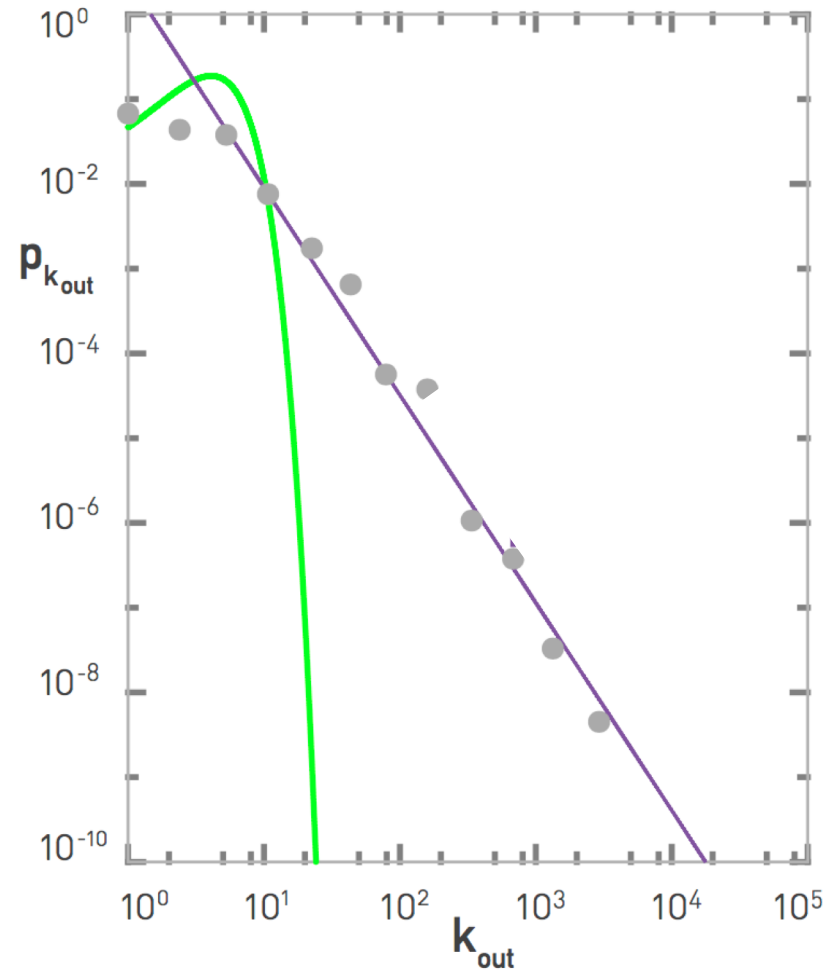
complementary cumulative  
distr. function (CCDF)



# Two degree distributions for directed networks



$p_{k_{in}}$  = fraction of nodes with  
input degree equal to  $k_{in}$



$p_{k_{out}}$  = fraction of nodes with  
output degree equal to  $k_{out}$



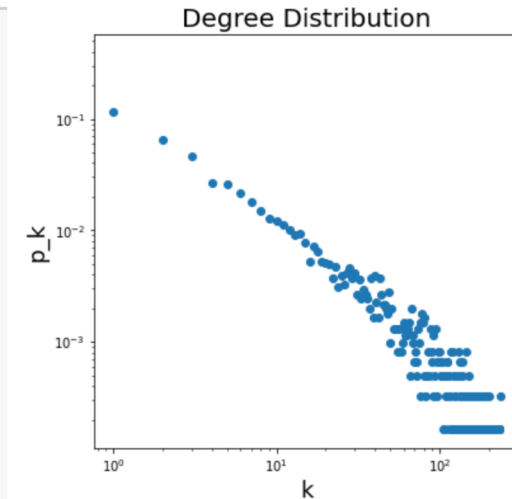


```
G = np.loadtxt('Wiki-Vote.txt').astype(int)

# adjacency matrix
N = np.max(G)
A = csr_matrix((np.ones(len(G)), (G[:, 1], G[:, 0])))

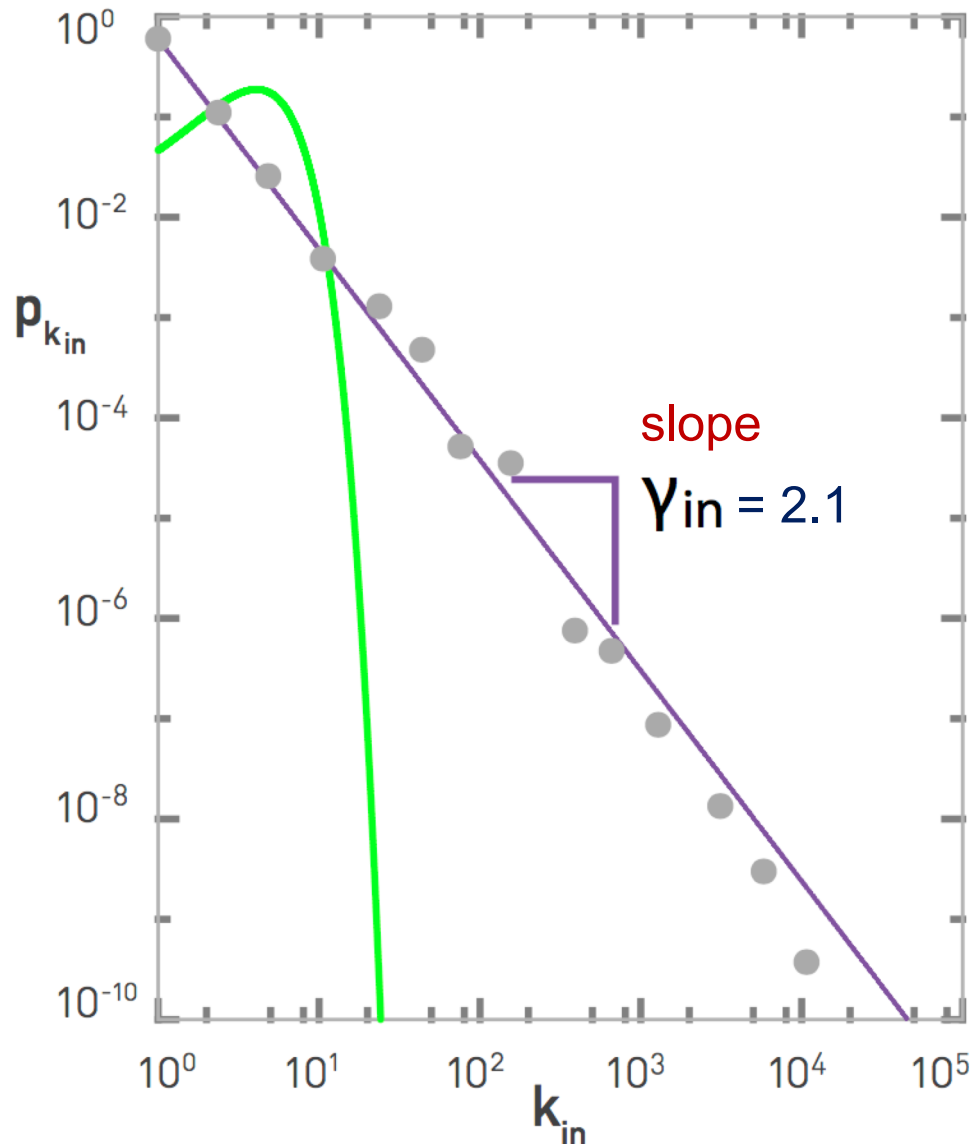
#distribution
which_deg = 0 # 0=out degree, 1=in degree
d = np.sum(A, which_deg) # out degree for each node
d = np.squeeze(np.asarray(d)) # from matrix to array
d = d[d>0] # avoid zero degree
k = np.unique(d) # degree samples
pk = np.histogram(d, k)[0] # occurrence of each degree
pk = pk/np.sum(pk) # normalize to 1
Pk = 1 - np.cumsum(pk) # complementary cumulative
```

```
fig = plt.figure()
plt.loglog(pk, 'o')
plt.title("Degree Distribution", size = 20)
plt.xlabel("k", size = 18)
plt.ylabel("p_k", size = 18)
plt.show()
```





# The power-law typical behaviour of social networks



many networks follow  
a **power-law**

$$\ln(p_k) = c - \gamma \cdot \ln(k)$$

$$p_k = C \cdot k^{-\gamma}$$

how to correctly  
estimate the **slope**  $\gamma$  ?



Degree distribution  $p_k = C k^{-\gamma}$

Constant  $C$  is determined by the (approx.)  
**normalization** condition

$$\int_{k_{\min}}^{\infty} p_k dk = C \cdot k_{\min}^{-(\gamma-1)} / (\gamma-1) = 1$$

Target PDF  $p(k|\gamma) = (\gamma-1)/k_{\min} \cdot (k/k_{\min})^{-\gamma}$



**ML criterion:** find the  $\gamma$  that best fits the data

$$\max_{\gamma} \sum_i \ln p(k_i | \gamma)$$

where  $k_i$  is the measured degree of node  $i$

$$f(\gamma) = \sum \ln((\gamma-1)/k_{\min}) - \gamma \ln(k_i/k_{\min})$$

$$f'(\gamma) = \sum 1/(\gamma-1) - \ln(k_i/k_{\min}) = 0$$

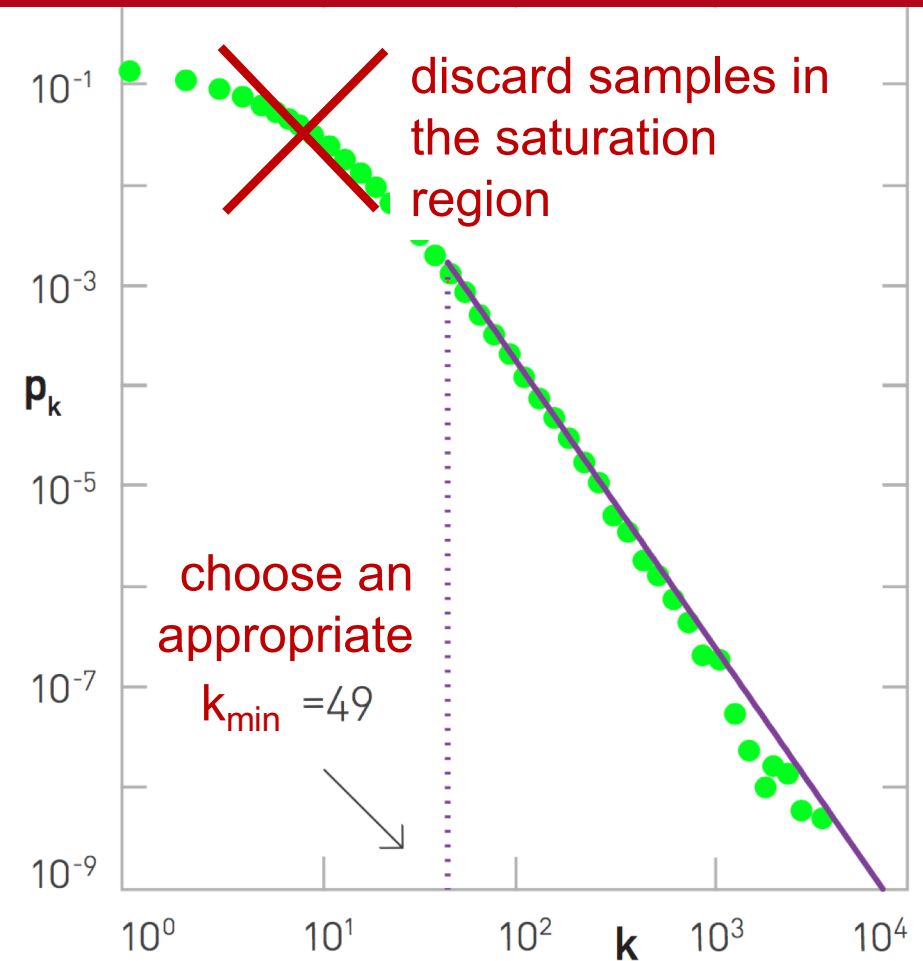
$$\gamma = 1 + \sum_i 1 / \sum_i \ln(k_i/k_{\min})$$



# Pseudocode example to estimate the exponent



```
which_deg = 1; % 1 = out degree  
d = full(sum(A,which_deg));  
d2 = d(d>=kmin); % restrict range  
ga = 1+1/mean(log(d2/kmin)); % estimate the exponent
```





# The value of the exponent $\gamma$ in real networks $\gamma \in [2,5]$

NETWORK	$N$	$L$	$\langle k \rangle$	$\gamma_{in}$	$\gamma_{out}$	$\gamma$
Internet	192,244	609,066	6.34	-	-	3.42*
WWW	325,729	1,497,134	4.60	2.00	2.31	-
Power Grid	4,941	6,594	2.67	-	-	Exp.
Mobile Phone Calls	36,595	91,826	2.51	4.69*	5.01*	-
Email	57,194	103,731	1.81	3.43*	2.03*	-
Science Collaboration	23,133	93,439	8.08	-	-	3.35*
Actor Network	702,388	29,397,908	83.71	-	-	2.12*
Citation Network	449,673	4,689,479	10.43	3.03**	4.00*	-
E. Coli Metabolism	1,039	5,802	5.58	2.43*	2.90*	-
Protein Interactions	2,018	2,930	2.90	-	-	2.89*

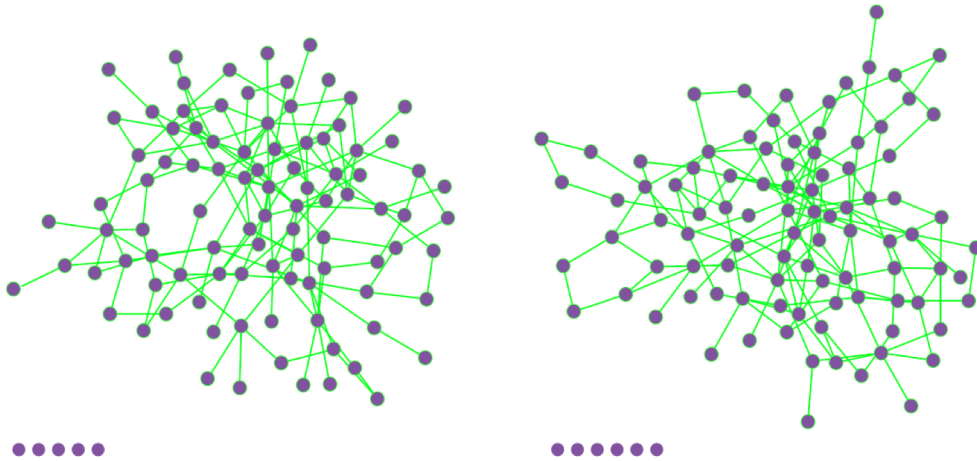
\* = good statistical fit with a **power-law**

\*\* = good fit for a power-law with an exponential cutoff

Exp = good fit with an exponential distribution  $e^{-ak}$

# Explaining the power-law

Preferential attachment



- ❑ The **random** network is the simplest model:
  - pick a probability  $p$ , with  $0 < p < 1$
  - activate each link  $(i,j)$  with probability  $p$
- ❑ The number of links is variable
- ❑ There might be **isolates**
- ❑ Easy to calculate fundamental parameters

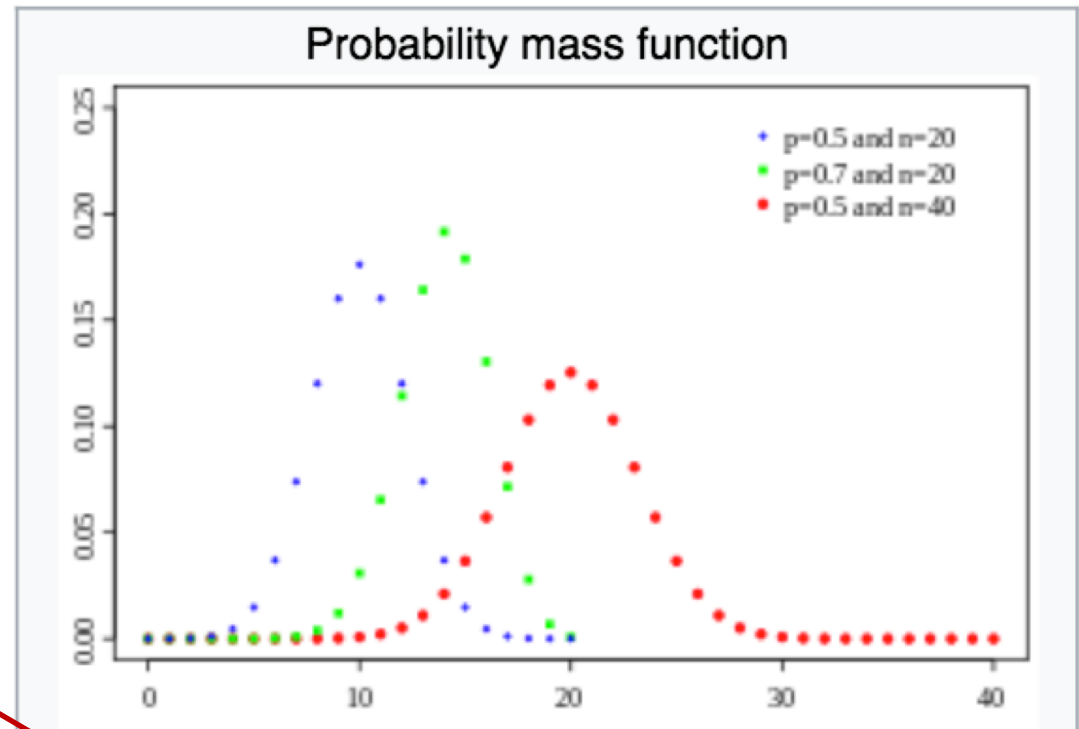




# Binomial distribution

explains the degree distribution for random networks

<b>Notation</b>	$B(n, p)$
<b>Parameters</b>	$n \in \{0, 1, 2, \dots\}$ – number of trials $p \in [0, 1]$ – success probability for each trial $q = 1 - p$
<b>Support</b>	$k \in \{0, 1, \dots, n\}$ – number of successes
<b>PMF</b>	$\binom{n}{k} p^k q^{n-k}$
<b>CDF</b>	$I_q(n - k, 1 + k)$
<b>Mean</b>	$np$
<b>Median</b>	$\lfloor np \rfloor$ or $\lceil np \rceil$
<b>Mode</b>	$\lfloor (n + 1)p \rfloor$ or $\lceil (n + 1)p \rceil - 1$
<b>Variance</b>	$npq$
<b>Skewness</b>	$\frac{q - p}{\sqrt{npq}}$
<b>Ex. kurtosis</b>	$\frac{1 - 6pq}{npq}$



$P(k;n,p)$  = probability that  $k$  out of  $n$  trials are positive, where each is positive with probability  $p$





- The number of neighbours is **binomially** distributed

$P(k;n,p)$  = probability that a node has **exactly**  $k$  neighbours, with number of possible neighbours  $n = N-1$

- Average # of neighbours

$$\langle k \rangle = (N-1)p \quad \rightarrow \quad p = \langle k \rangle / (N-1)$$

this defines  $p$

- Variance

$$\sigma_x^2 = (N-1)p(1-p) \simeq \langle k \rangle$$

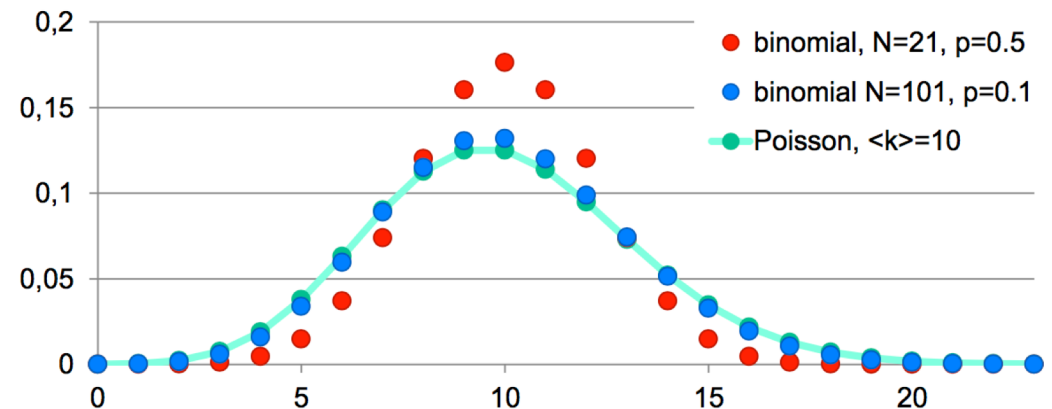
$p$  is usually very small (since  $\langle k \rangle \ll N$ )

tight around the mean



## □ Poisson distribution (easier to use)

$$P[x = k] = \frac{m_x^k}{k!} \cdot e^{-m_x}$$

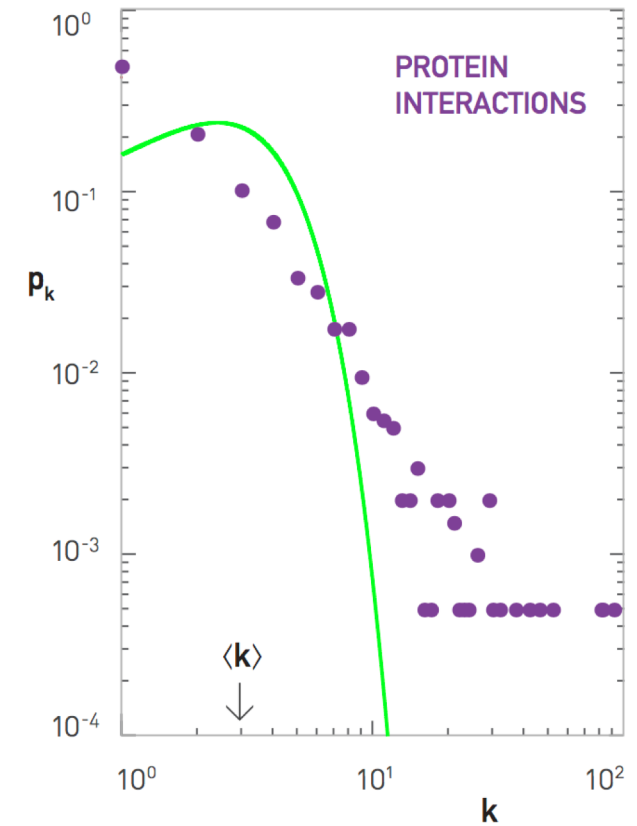
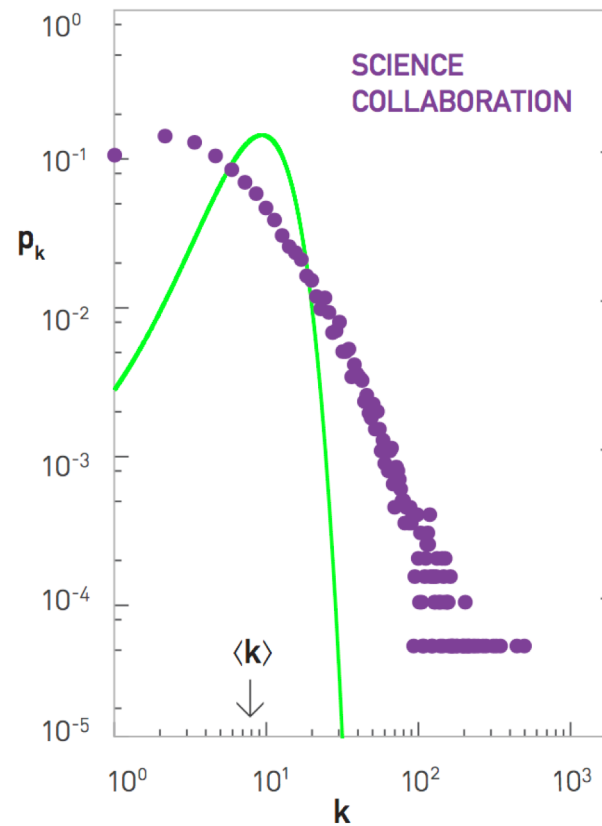
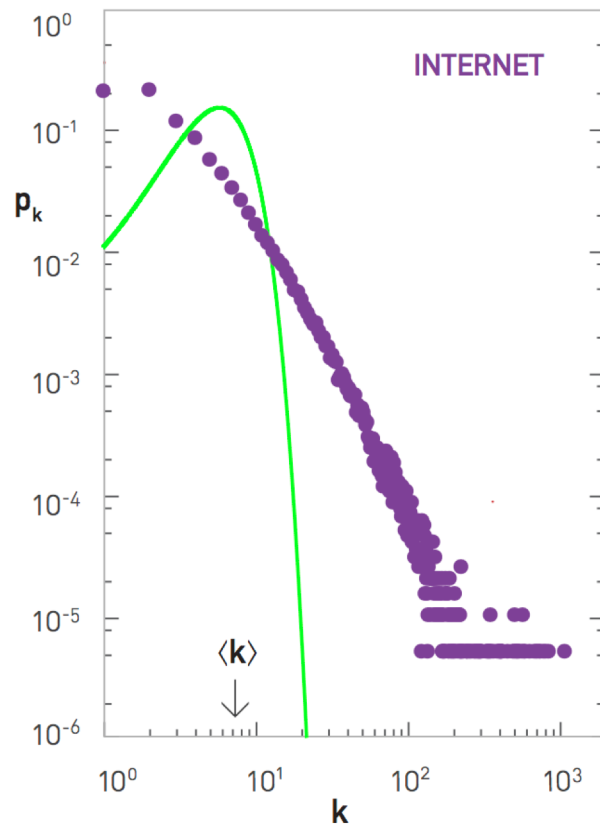


## □ Very good approximation of binomial for small $p$ (and at small $k$ )

$$P[x = k] = \underbrace{\frac{(n - k + 1) \dots (n - 1)n}{n^k}}_{\simeq 1} \cdot \overset{\text{active part}}{\downarrow} \frac{m_x^k}{k!} \cdot \underbrace{\left(1 - \frac{m_x}{n}\right)^{n-k}}_{\simeq \text{const}} \overset{p}{\downarrow}$$



# Are real networks Poisson? no, they aren't

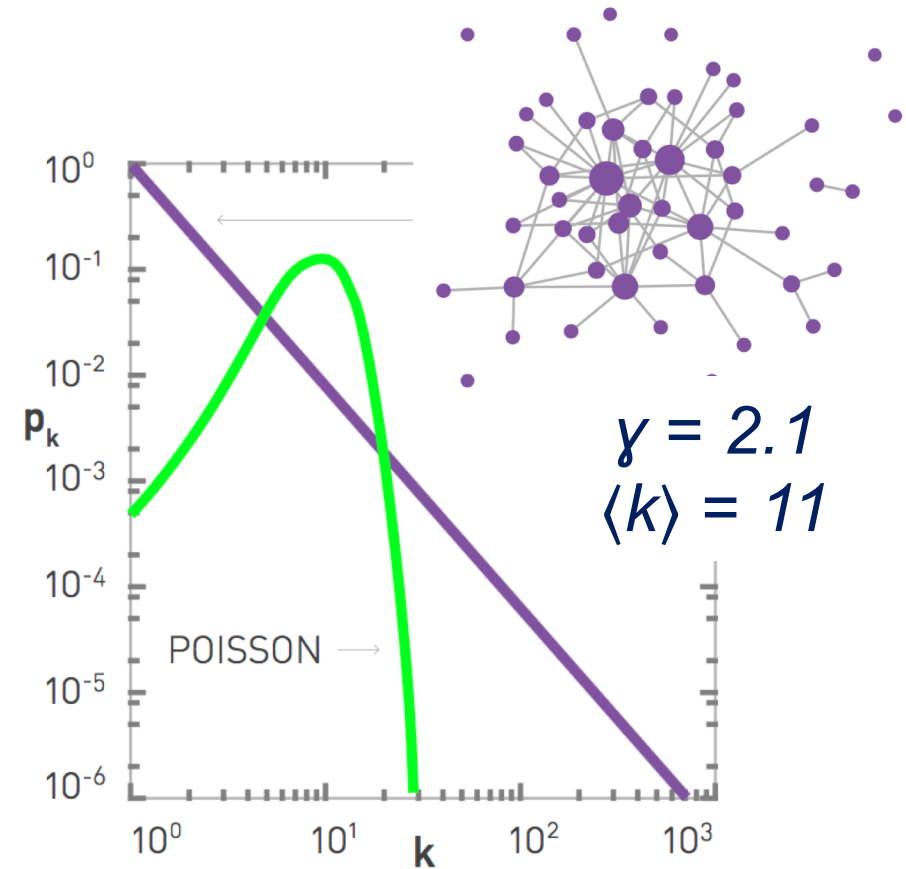
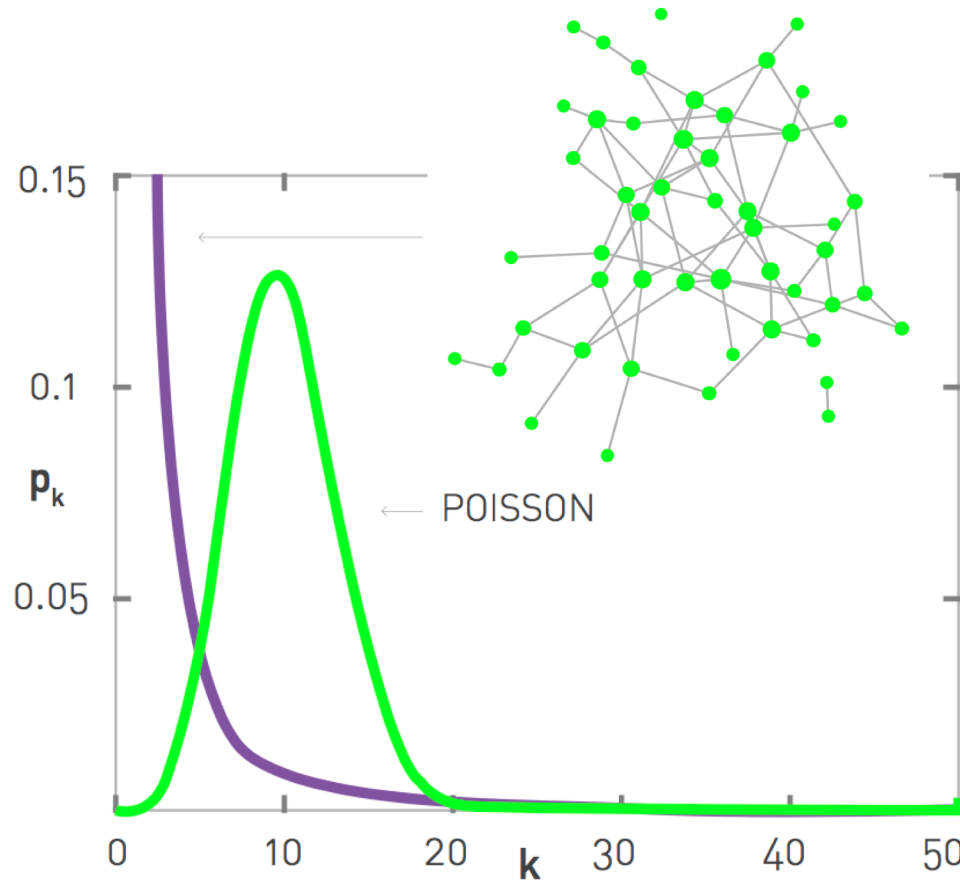


**No!** Poisson networks are deprived of hubs

... but, nevertheless, Poisson networks capture some aspects



# Poisson versus power law a comparison



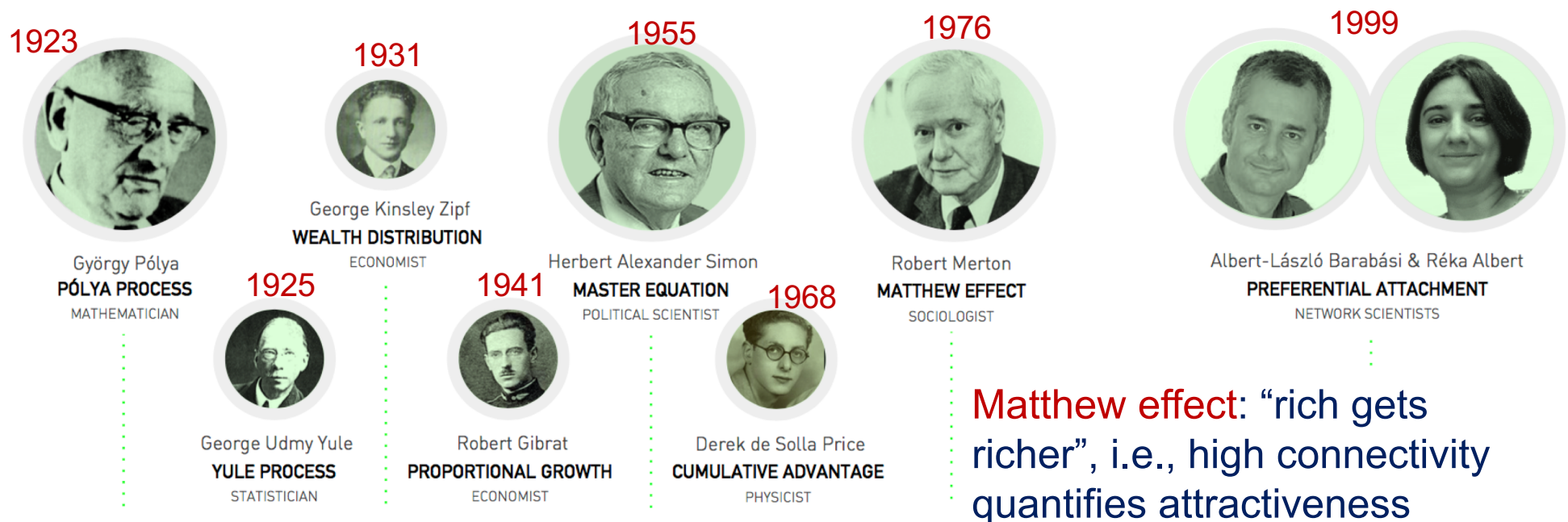
Power-law is **heavy tailed** (presence of hubs) -  
like Weibull, lognormal, Lévy



## Nodes link to the **more connected** nodes

e.g., think of www

## This idea has a long history





Start with  $m_0$  nodes arbitrarily connected, with  $\langle k \rangle = m$

## □ Growth

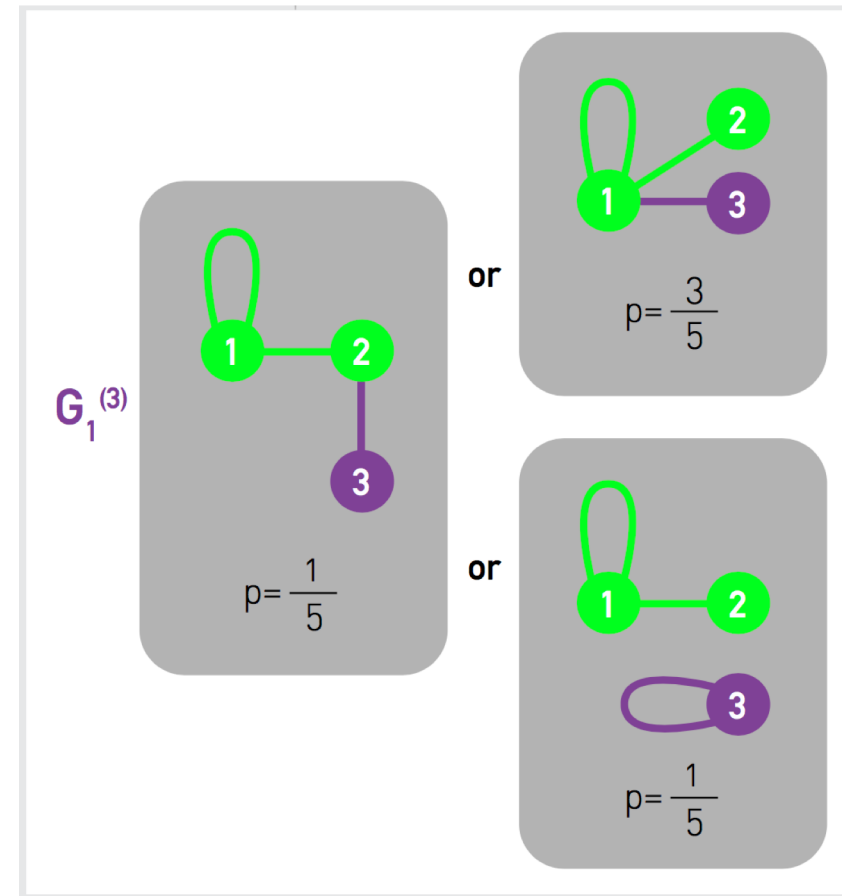
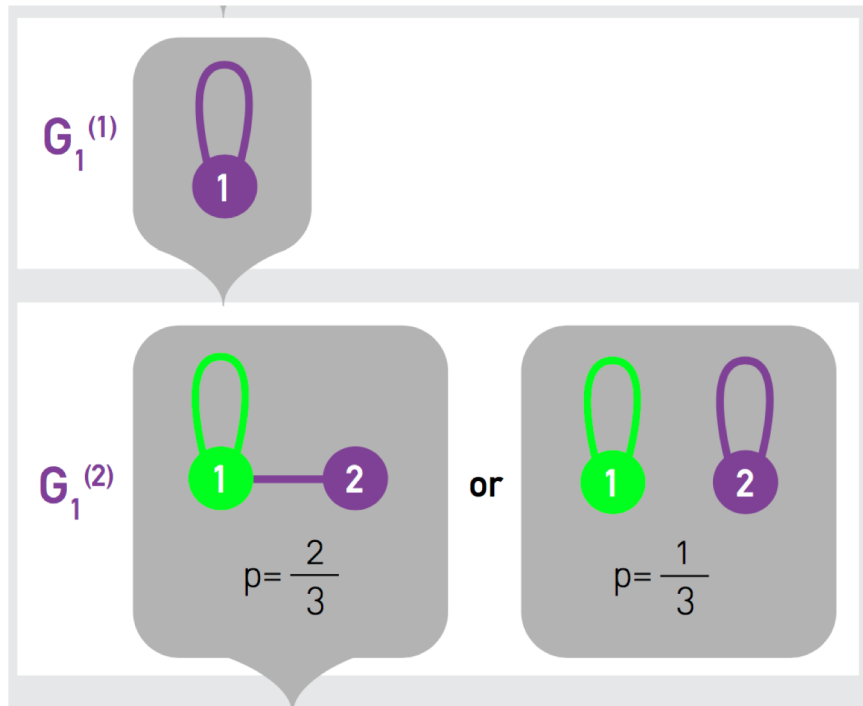
add a node (the  $N$ th) with  $m$  links that connect the node to nodes in the network

## □ Preferential attachment

$p_i = k_i / C$  probability of connecting to node  $i$

$p_i = 1/C$  for self-loops

$$C = 1 + \sum k_i = 1 + 2(N-1)m$$







- Increase in the degree (at each step)

$$\Delta k_i \simeq \underset{\substack{\uparrow \\ \text{trials}}}{m} \cdot \underset{\substack{\uparrow \\ \text{probability per trial}}}{k_i} / (1+2m(N-1)) \simeq k_i / 2N$$

- Approximation in the continuous domain

$$\Delta k_i \simeq dk_i/dN \rightarrow dk_i/k_i \simeq 1/2 dN/N$$

- Integration

$$\ln(k_i) = 1/2 \ln(N) + \text{const.} \rightarrow k_i = c N^{1/2}$$

- Recalling that node  $i$  joins the network at time  $N = i$

$$k_i(N=i) = m \rightarrow k_i(N) = m (N/i)^{1/2}$$

1/2 is the  
dynamic  
exponent



- Recall  $k_i = m (N/i)^{1/2}$
- The number of nodes with degree smaller than  $k$  is

$$k_i < k \rightarrow m (N/i)^{1/2} < k$$

$$\rightarrow i > N (m/k)^2 \rightarrow N - N (m/k)^2$$

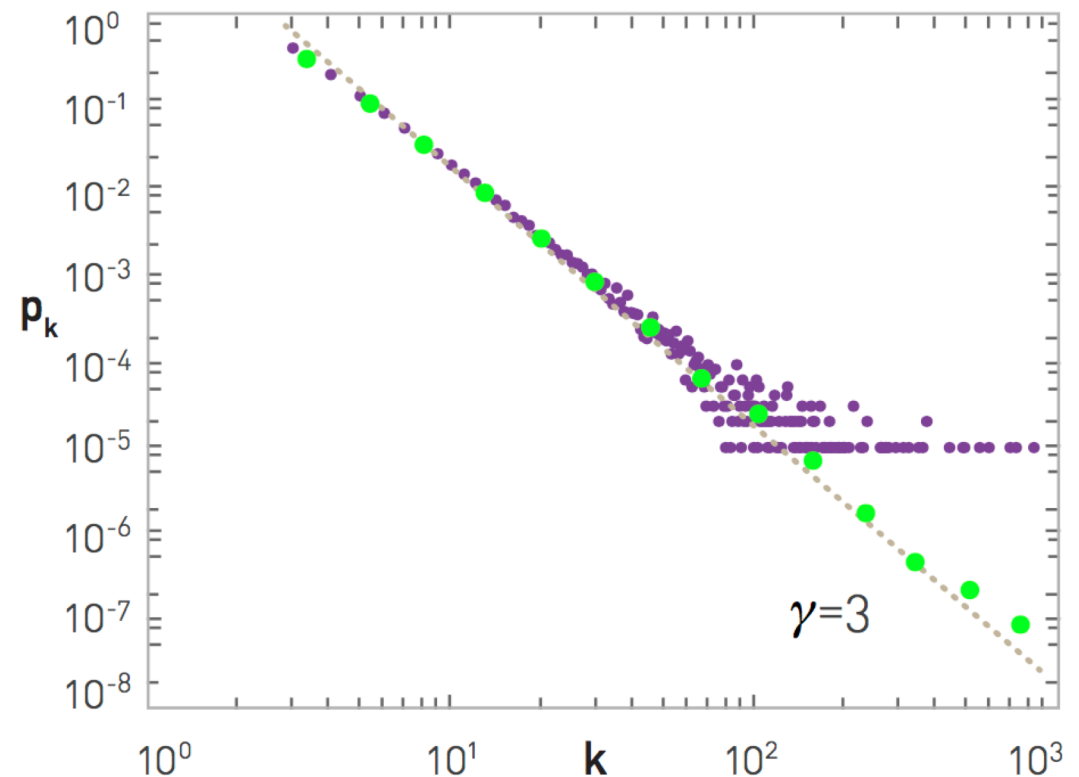
- CDF is  $P_k = P[k_i \leq k] = 1 - (m/k)^2$

- The degree distribution is

$$dP_k / dk = p_k = 2 m^2 / k^3$$



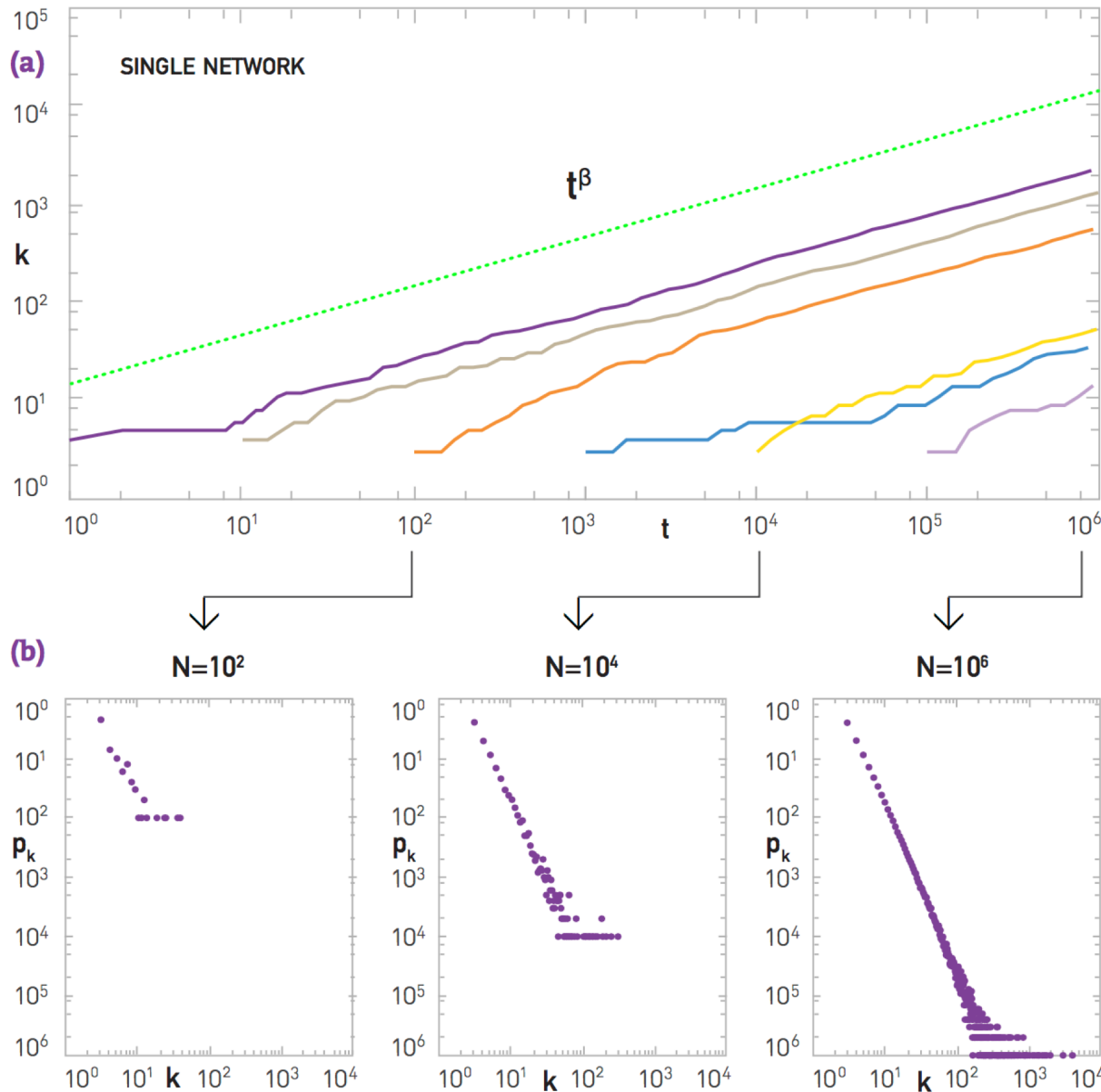
- Depending on the implementation there might be **self/multiple** links
- Most nodes have a small degree (exactly  $m$  for the youngest ones)
- Hubs appear
- The average degree is  $\langle k \rangle = 2m$ , and in fact  $L = Nm = \frac{1}{2}\langle k \rangle N$
- The resulting degree distribution is always a power-law with exponent  $\gamma = 3$





# The Barabasi-Albert model

consequence of  $k_i = m (N/i)^{1/2}$



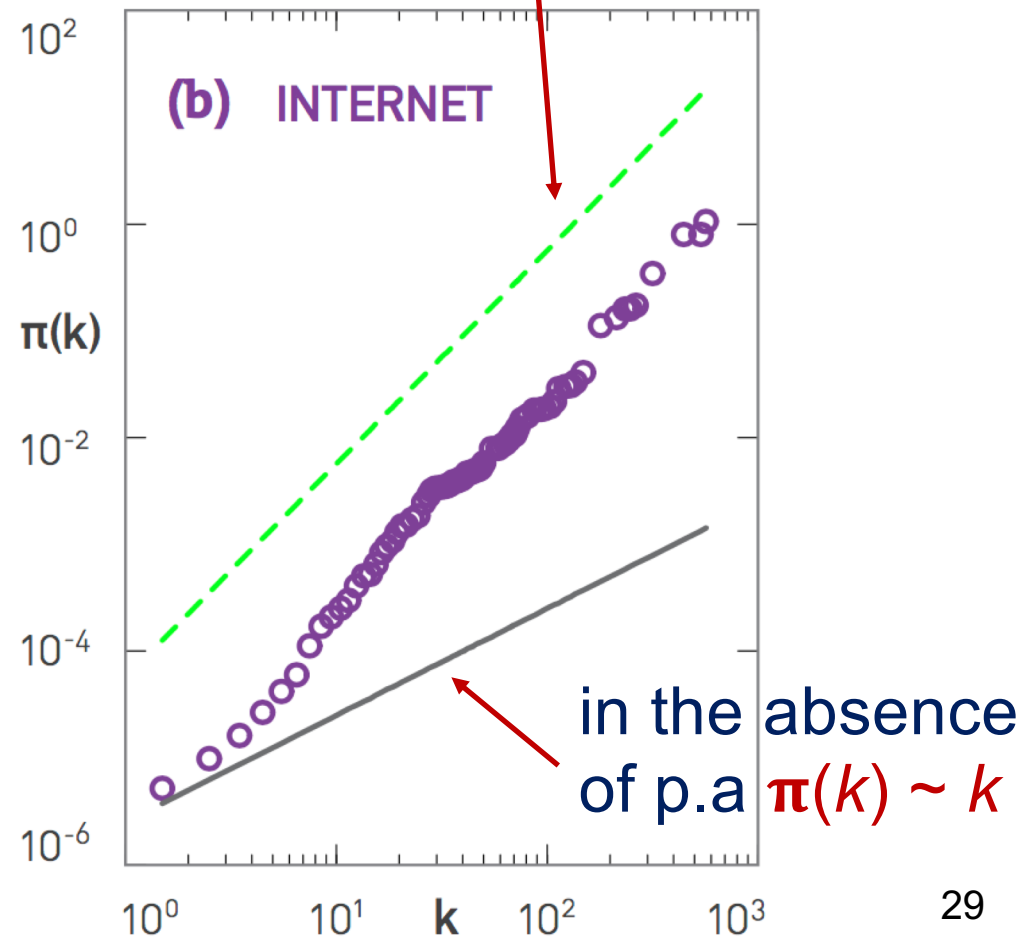
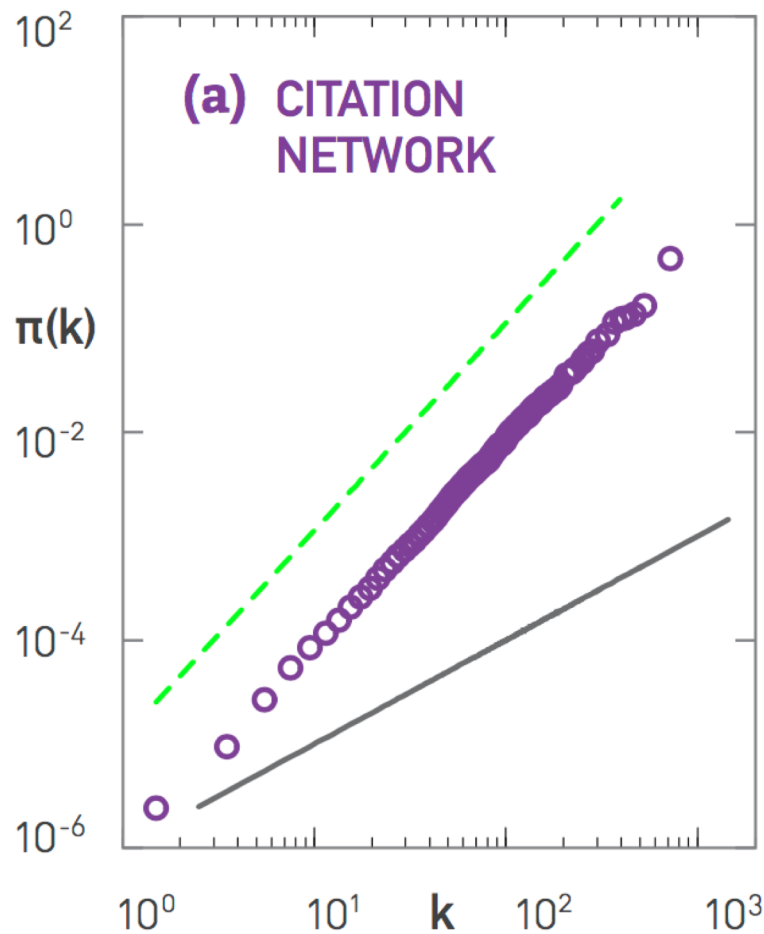
- ❑ all nodes follow the same dynamics
- ❑ the growth is **sub-linear**: nodes are competing with the others
- ❑ the earlier the node is added, the higher the degree – “**first-mover advantage**”
- ❑ older nodes acquire more links
- ❑ this explains the hub formation



# Measuring preferential attachment in real networks

$$\pi(k) = \sum_{k_i}^k \Delta k_i / \Delta N$$

under preferential attachment  $\pi(k) \sim k^2$





# The Bianconi-Barabasi model

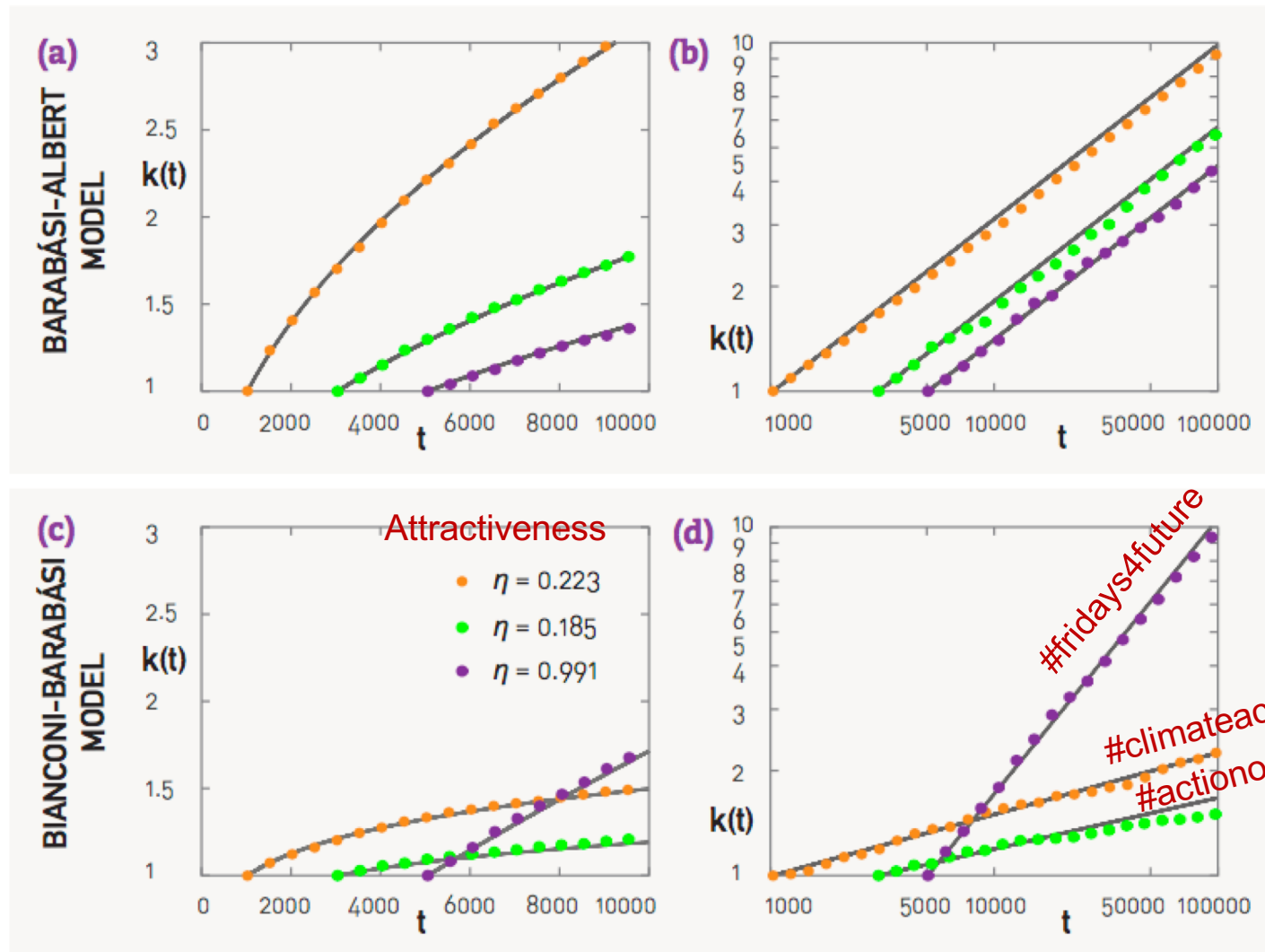
Bianconi, Barabási. "Competition and multiscaling in evolving networks" (2001)

The model:

- ❑ **Growth** – at time step  $N$  a new node  $i=N$  is added with  $m$  **links** and **fitness**  $\eta_i$
- ❑ Attractiveness (or fitness) is a random number drawn from a given **distribution**  $\rho(\eta)$  - a quality of the individual to attract links
- ❑ **Preferential attachment** - probability of linking to node  $i$  is proportional to both the degree and the attractiveness, i.e.,  $p_i = k_i \eta_i / \sum k_j \eta_j$



we guess  $k_i \simeq m (N/i)^{\beta(\eta)}$  for some  $\beta(\eta)$





□ We guess  $k_i \simeq m (N/i)^{\beta(\eta_i)}$

□ Increase in the degree  $\Delta k_i \simeq \overset{\text{trials}}{m} \cdot \overset{\text{probability per trial}}{k_i \eta_i} / \sum k_j \eta_j$

□ We show that  $\sum k_j \eta_j \simeq m N \cdot C$  (see proof)





- Analysis of denominator  $\sum k_i \eta_i$

→ average value wrt  $\eta$

→ hypothesis  $k_i \simeq m (N/i)^{\beta(\eta)}$

- $A = E[ \sum_i k_i \eta_i ] = \sum E[ k_i \eta_i ] \simeq \int_1^N E[k_i \eta_i] di$

- $E[k_i \eta_i] = \int m(N/i)^{\beta(\eta)} \eta \cdot \rho(\eta) d\eta$

- Swap integrals

$$A \simeq \int m N^{\beta(\eta)} \left[ \int_1^N i^{-\beta(\eta)} di \right] \eta \cdot \rho(\eta) d\eta$$

- Integrate

$$A \simeq m N \cdot \int \frac{(1 - N^{\beta(\eta)-1}) \eta \rho(\eta) d\eta}{1 - \beta(\eta)}$$

constant C

negligible for large  $N$  if  $0 < \beta < 1$



- We guess  $k_i \simeq m (N/i)^{\beta(\eta_i)}$
- Increase in the degree  $\Delta k_i \simeq m \cdot k_i \eta_i / \sum k_j \eta_j$
- It is  $\sum k_j \eta_j \simeq m N \cdot C$

Hence:

1. By inspection of the above

$$\Delta k_i \simeq m (N/i)^{\beta(\eta_i)} \eta_i / N C$$

2. By continuum theory

$$\Delta k_i \simeq dk_i/dN \simeq m \beta(\eta_i) N^{\beta(\eta_i) - 1} i^{-\beta(\eta_i)}$$

3. By combining the results  $\beta(\eta_i) \simeq \eta_i/C$

We conclude  $k_i \simeq m (N/i)^{\eta_i/C}$



$$\beta(\eta) \simeq \eta / C$$

$$C = \frac{\int \eta \rho(\eta) d\eta}{1 - \beta(\eta)} \rightarrow$$

$$1 = \int_0^{\eta_{\max}} (C - \eta)^{-1} \eta \rho(\eta) d\eta$$

this identifies C for a given  $\rho(\eta)$

it is  $C > \eta_{\max}$ , i.e.,  $\beta < 1$ ,  $\rightarrow$  the integral makes sense

growth with  
exponent  $< 1$

it also is  $C \leq 2\eta_{\max}$



Want to identify  $P_k = P[k_i \leq k] = 1 - P[k_i > k]$

- $k_i > k$  and  $k_i = m (N/i)^{\eta/C} \rightarrow i < N (m/k)^{C/\eta}$
- Hence  $P[k_i > k | \eta_i] = (m/k)^{C/\eta_i}$
- and  $P[k_i \leq k | \eta_i] = 1 - (m/k)^{C/\eta_i}$
- We have  $P_k = 1 - \int (m/k)^{C/\eta} \rho(\eta) d\eta$

The degree distribution is

$$p_k = P_k' = C \int_0^{\eta_{\max}} k^{-(C/\eta+1)} m^{C/\eta} \eta^{-1} \rho(\eta) d\eta$$

weighted combination of power laws with  
exponent in  $[2, \infty)$  since  $\eta_{\max} < C$



What if  $\rho(\eta) = \delta(\eta-1)$  ?

□ Coefficient  $C = 2$  since

$$\int_0^{\eta_{\max}} (C/\eta - 1)^{-1} \delta(\eta-1) d\eta = (C - 1)^{-1} = 1$$

□ Exponential degree  $k_i \simeq m (N/i)^{1/2}$

Degree distribution

$$p_k = C \int_0^{\eta_{\max}} \eta^{-1} m^{C/\eta} k^{-(C/\eta+1)} \delta(\eta-1) d\eta = 2 m^2 k^{-3}$$



What if  $\rho(\eta) = 1$  and  $\eta_{\max} = 1$  ?

□ Coefficient  $C = 1.255$  since

$$\int_0^1 (C/\eta - 1)^{-1} d\eta = 1 \rightarrow e^{-2/C} = 1 - 1/C$$

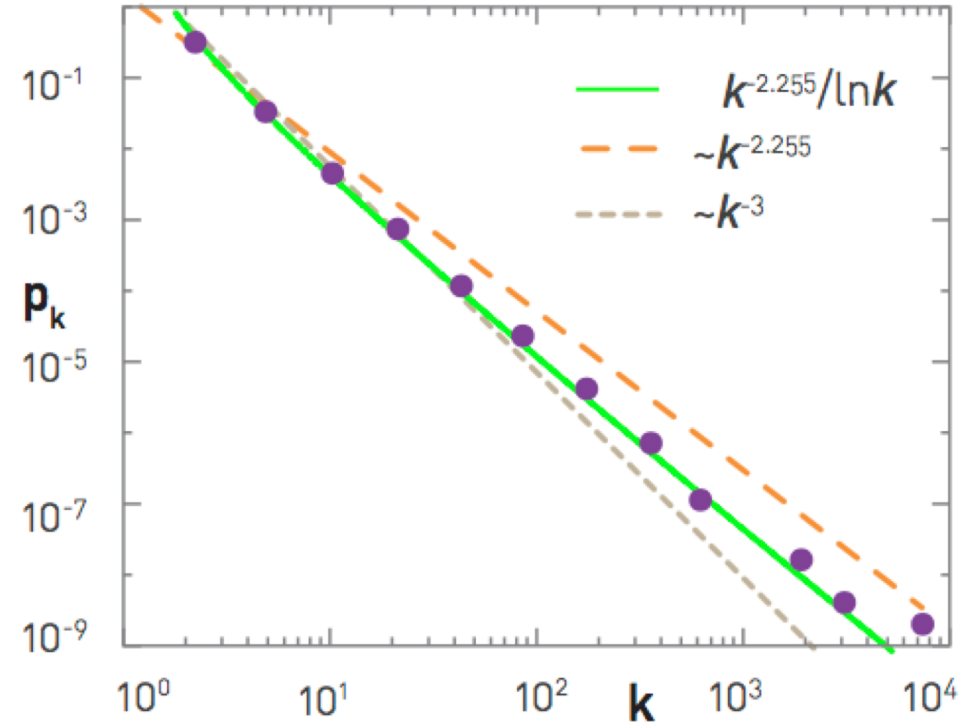
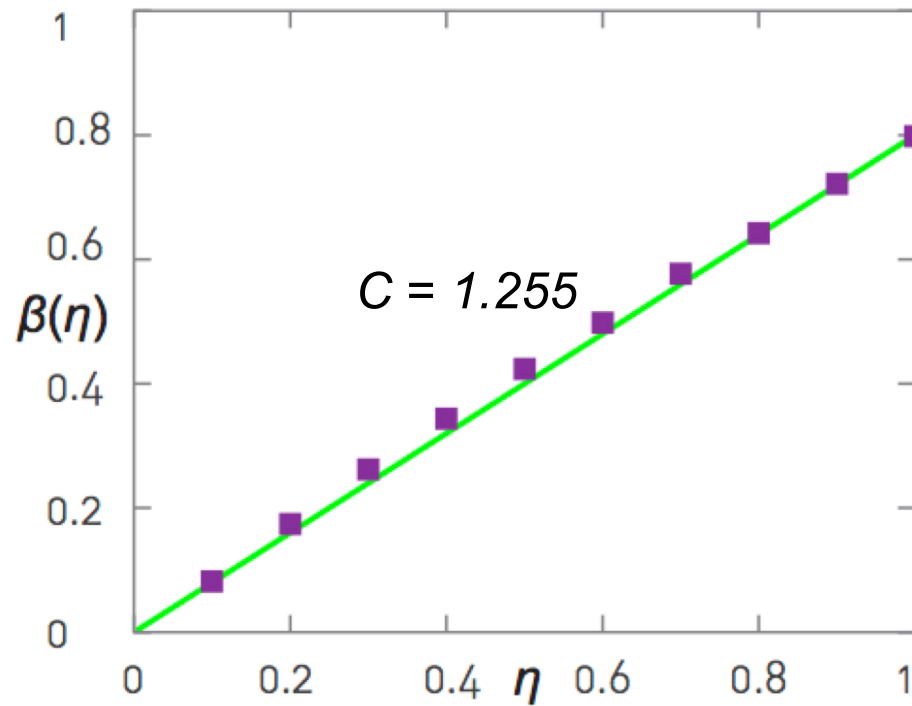
□ Exponential degree  $k_i \simeq m (N/i)^{\eta_i/C}$

□ Each node has its **own dynamic exponent !!!**

Degree distribution

$$p_k = C/k \int_0^1 \eta^{-1} e^{-C \ln(k/m)/\eta} d\eta \sim k^{-(1+C)} / \ln(k)$$

$e^{-b} - b E_1(b)$ ,  $b = C \ln(k/m)$   
exponential integral  $E_1$



degree distribution  $p_k \sim k^{-2.255} / \ln(k)$

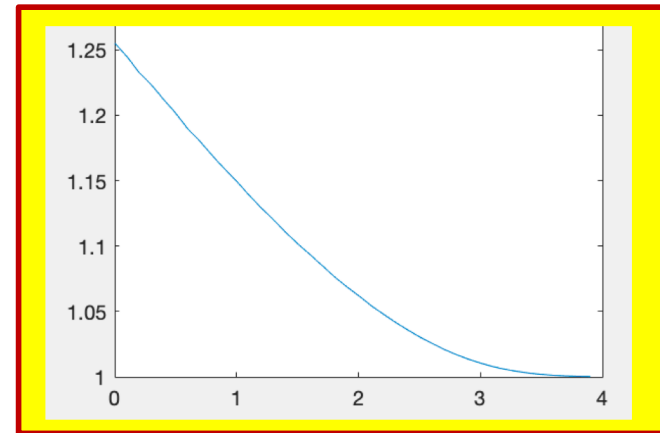
corrective term



What if  $\rho(\eta) = a e^{-a\eta} / (1-e^{-a})$  and  $\eta_{\max} = 1$  ?

- C rapidly converges to C=1

$$\int_0^1 (C/\eta - 1)^{-1} \rho(\eta) d\eta = 1$$



- Exponential degree  $k_i \simeq m (N/i)^{\eta/C}$
- Each node has its own dynamic exponent !!!

Degree distribution

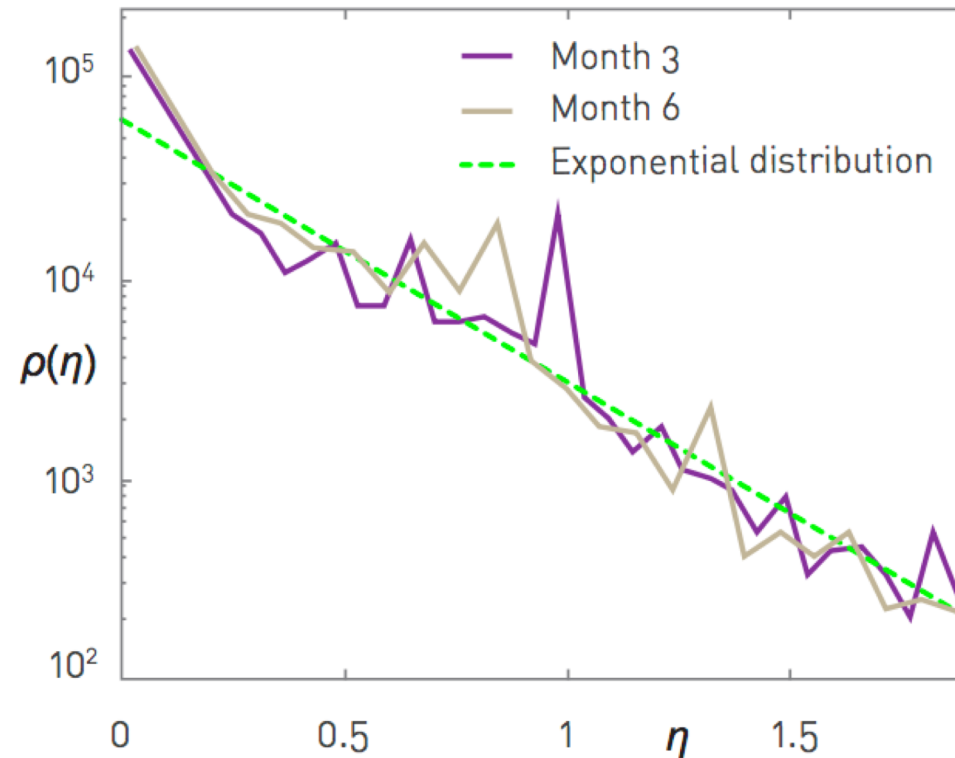
$$p_k = C/k \int_0^1 \eta^{-1} e^{-C \ln(k/m)/\eta} \rho(\eta) d\eta \sim k^{-(1+C)} / \ln(k)$$

↑  
exponential integral E<sub>1</sub>





$$a = 4.6$$
$$C = 1$$

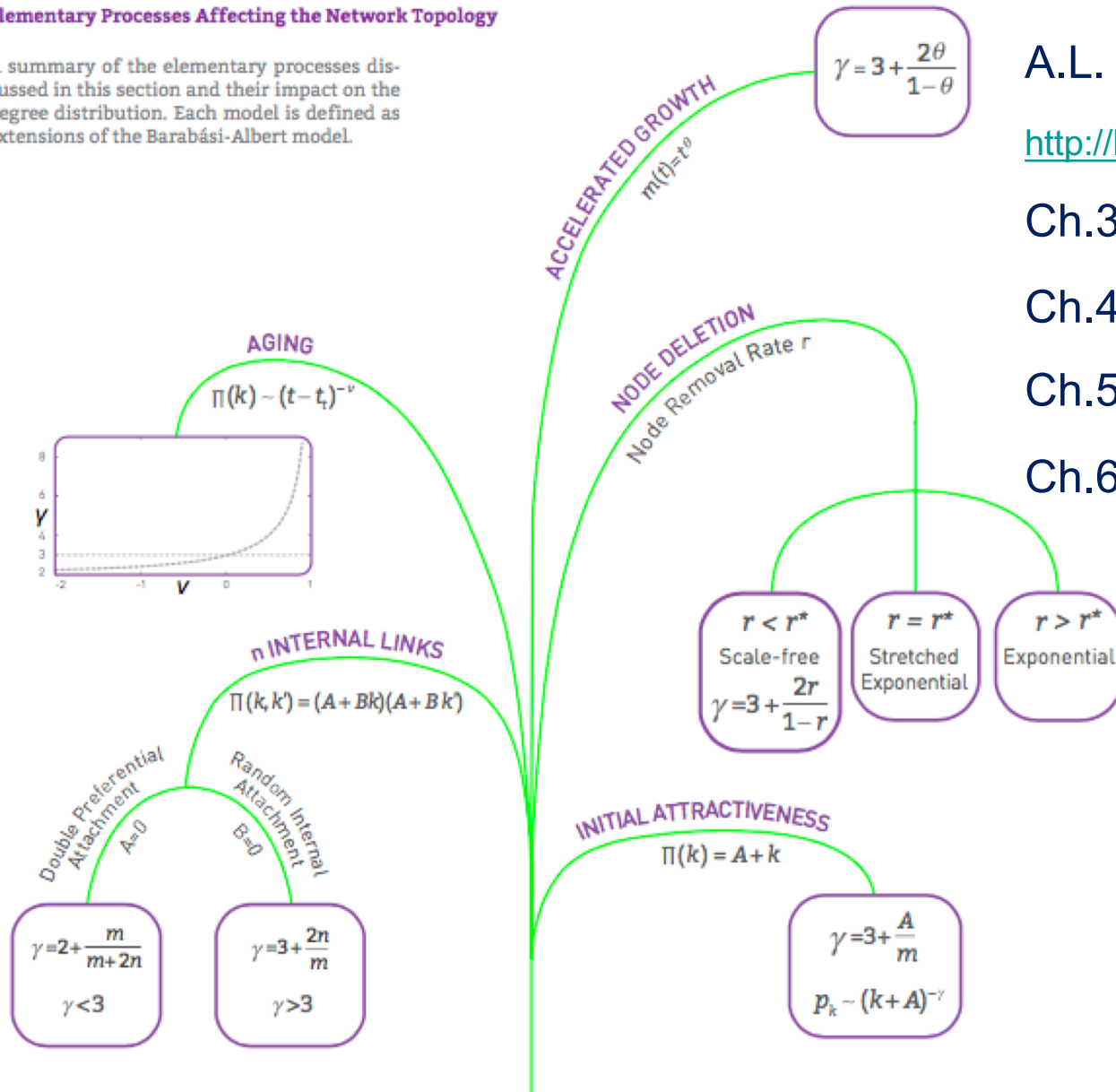


degree distribution  $p_k \sim k^{-2} / \ln(k)$



## Elementary Processes Affecting the Network Topology

A summary of the elementary processes discussed in this section and their impact on the degree distribution. Each model is defined as extensions of the Barabási-Albert model.



A.L. Barabási, Network science

<http://barabasi.com/networksciencebook>

Ch.3 “Random networks”

Ch.4 “The scale-free property”

Ch.5 “The Barabási-Albert model”

Ch.6 “Evolving networks”

# Properties of the power-law

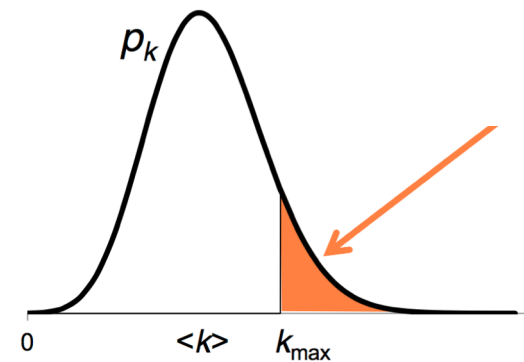
scale-free and random networks



Degree distribution  $p_k = C k^{-\gamma}$  with  $C = (\gamma-1) k_{\min}^{\gamma-1}$

The size of the largest hub is captured by

$$\int_{k_{\max}}^{\infty} p_k dk = C \cdot k_{\max}^{-(\gamma-1)} / (\gamma-1) = 1/N$$



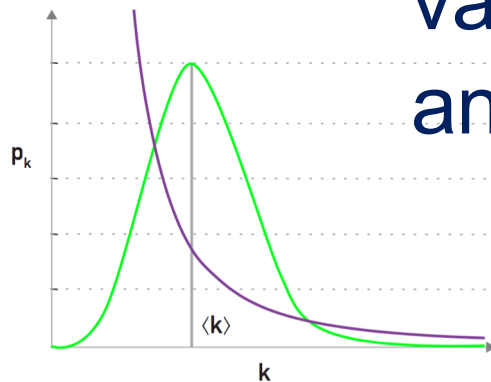
$$k_{\max} = k_{\min} N^{1/(\gamma-1)}$$

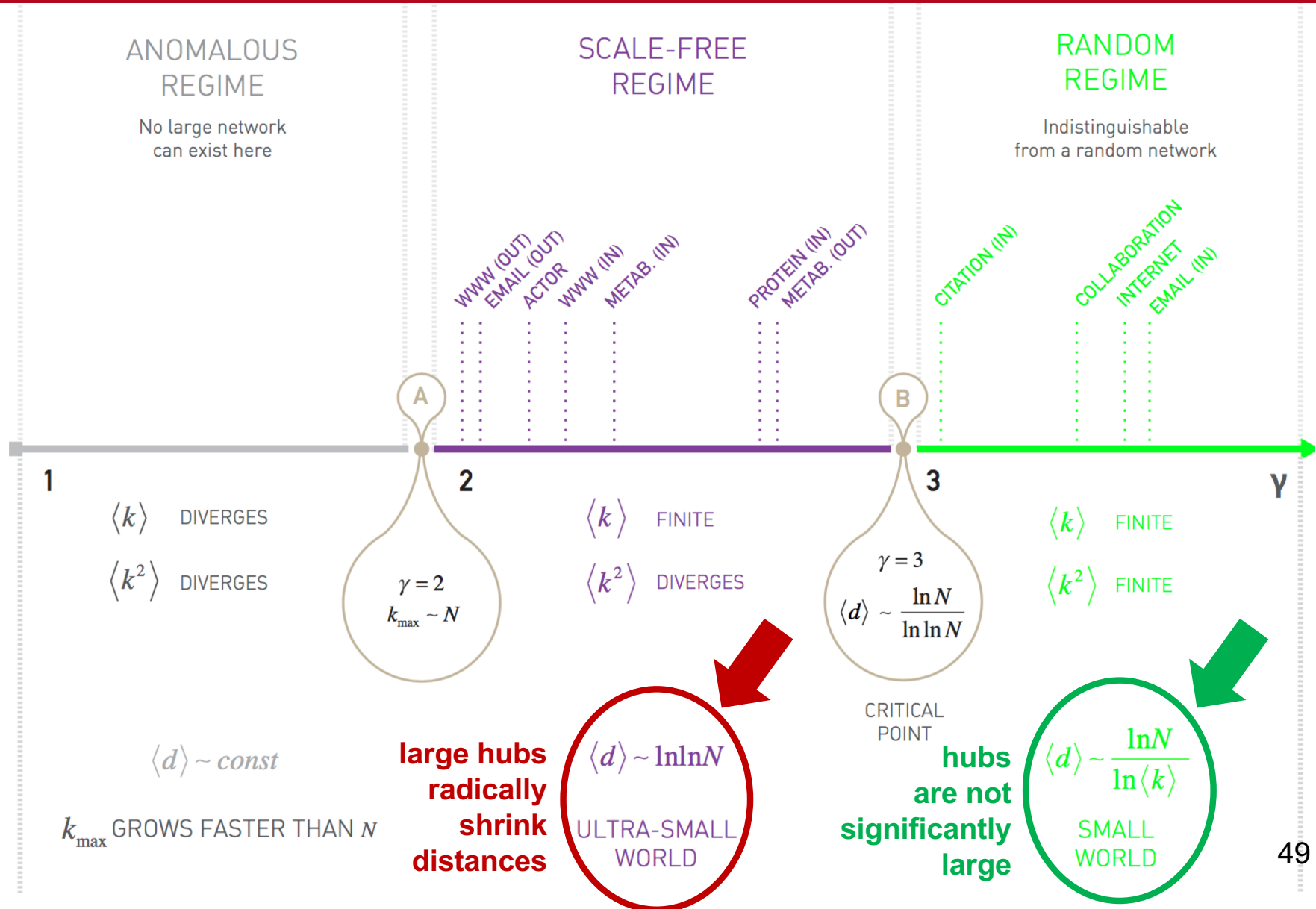
is the **natural cutoff**  
it explains large hubs



$$\begin{aligned} \square \langle k^n \rangle &= \int_{k_{\min}}^{k_{\max}} k^n p_k dk && \text{with } p_k = C k^{-\gamma} \\ &= C \cdot (k_{\max}^{n-\gamma+1} - k_{\min}^{n-\gamma+1}) / (n-\gamma+1) \\ &= C k_{\min}^{n-\gamma+1} \cdot (N^{-1+n/(\gamma-1)} - 1) / (n-\gamma+1) \end{aligned}$$

- $\square$  They diverge with  $N$  if  $\gamma < n+1$   
mean ( $n=1$ ) doesn't diverge for  $\gamma \geq 2$   
variance ( $n=2$ ) diverges for  $\gamma < 3$   
and the network does not have a scale  
(scale-free regime)





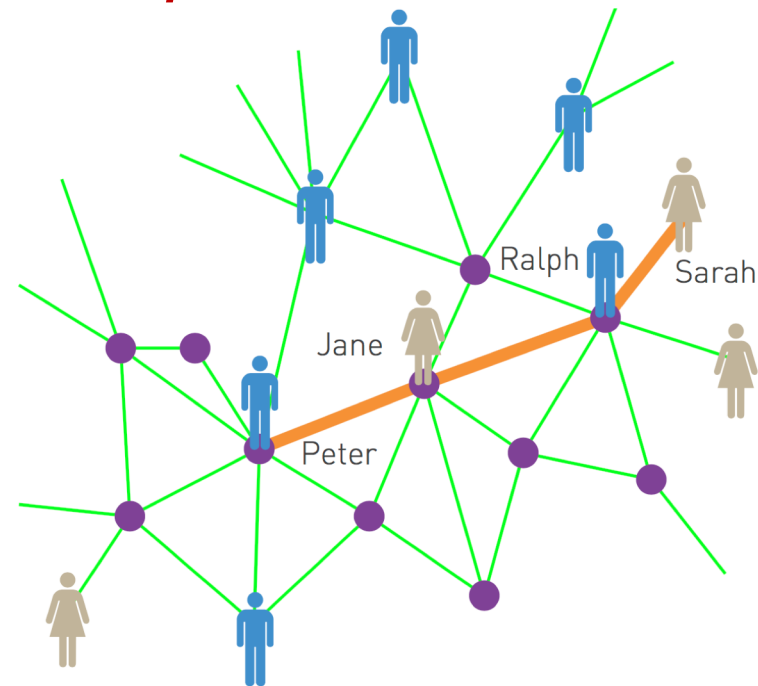


# Small world property

Watts, Strogatz, «*Collective dynamics of small-world networks*», (1998)

In real networks distance between two randomly chosen nodes is generally short

Milgram [1967]: *6 degrees of separation*



What does this mean?

We are more connected than we think



- we reach  $\langle k \rangle$  nodes in one hop,  $\langle k \rangle^2$  in two,  $\langle k \rangle^3$  in three, etc.
- an **estimate** of the average distance  $\langle d \rangle$  is found by solving for  $N = \langle k \rangle^{\langle d \rangle}$  to have

$$\langle d \rangle = \ln(N) / \ln(\langle k \rangle)$$

- $\langle d \rangle$  is often taken as an estimate of the network **diameter**  $d_{\max}$

e.g.: on earth we are  $N=7 \cdot 10^9$  individuals,

with  $\langle k \rangle=1000$  acquaintances each  $\rightarrow \langle d \rangle = 3.28$





# Distances in random graphs

fitting with real data

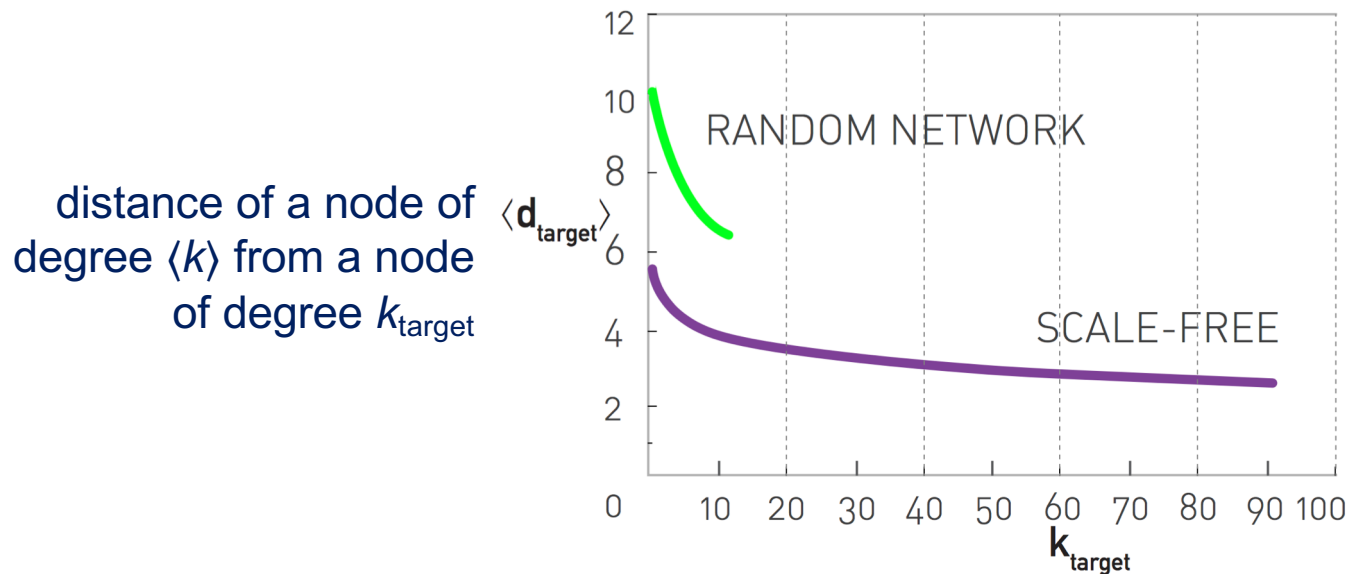
NETWORK	$N$	$L$	$\langle k \rangle$	$\langle d \rangle$	$d_{max}$	$\frac{\ln N}{\ln \langle k \rangle}$
Internet	192,244	609,066	6.34	6.98	26	6.58 ✓
WWW	325,729	1,497,134	4.60	11.27	93	8.31 ✓
Power Grid	4,941	6,594	2.67	18.99	46	8.66
Mobile Phone Calls	36,595	91,826	2.51	11.72	39	11.42 ✓
Email	57,194	103,731	1.81	5.88	18	18.4
Science Collaboration	23,133	93,439	8.08	5.35	15	4.81 ✓
Actor Network	702,388	29,397,908	83.71	3.91	14	3.04 ✓
Citation Network	449,673	4,707,958	10.43	11.21	42	5.55
E. Coli Metabolism	1,039	5,802	5.58	2.98	8	4.04
Protein Interactions	2,018	2,930	2.90	5.61	14	7.14 ✓

Very good fit ! Correct at least as order of magnitude



- The average distance increases as  $\ln(\ln(N))$ , much slower than  $N$  or  $\ln(N)$

e.g. in www  $N=7 \cdot 10^9$ ,  $\ln(N)=22.7$ ,  $\ln(\ln(N))=3.12$  (very small)



- The **large** hubs radically shrink the distance between nodes  $\rightarrow$  ultra small world



In many social experiments people avoided hubs for entirely perceptual reasons (e.g., they assumed they are busy, better use them only if really needed)

We live in a **ultra-small-world**, but we perceive that we are more distant from others than we really are!



# Friendship paradox

my friends are more popular than me (Feld 1991)

- ❑ Can be observed in the **ultra-small-world** under the presence of big hubs
- ❑ Rationale: a node is very likely to be connected to a big hub, having a very large number of connections
- ❑ # of friends (in the average) =  $\langle k \rangle$
- ❑ # of friends of friends  $\approx N$



- ❑ Do not use it for resizing nodes according to their importance (will use PageRank for this)
- ❑ Provide useful information in the form of a **degree distribution**
- ❑ Always plot degree distributions in the **log scale**
- ❑ Always evaluate their **slope  $\gamma$** , but please use the ML approach:  $\gamma$  provides useful insights on the network
- ❑ **Preferential attachment** and **attractiveness** can be measured if you have temporal info on the network

# PageRank centrality

Google's approach to centrality



# How to organise the web?

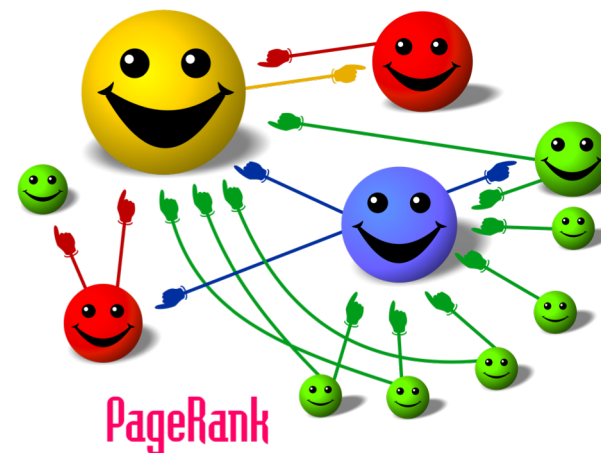
links as votes

- ❑ the higher (and stronger) the **number of incoming links**, the more important a node
- ❑ the more important a node, the more **valuable** the output links





- ❑ PageRank works by counting the number and quality of links to a page to determine a rough estimate of how important the website is
- ❑ The underlying assumption is that more important websites are likely to receive more links from other websites



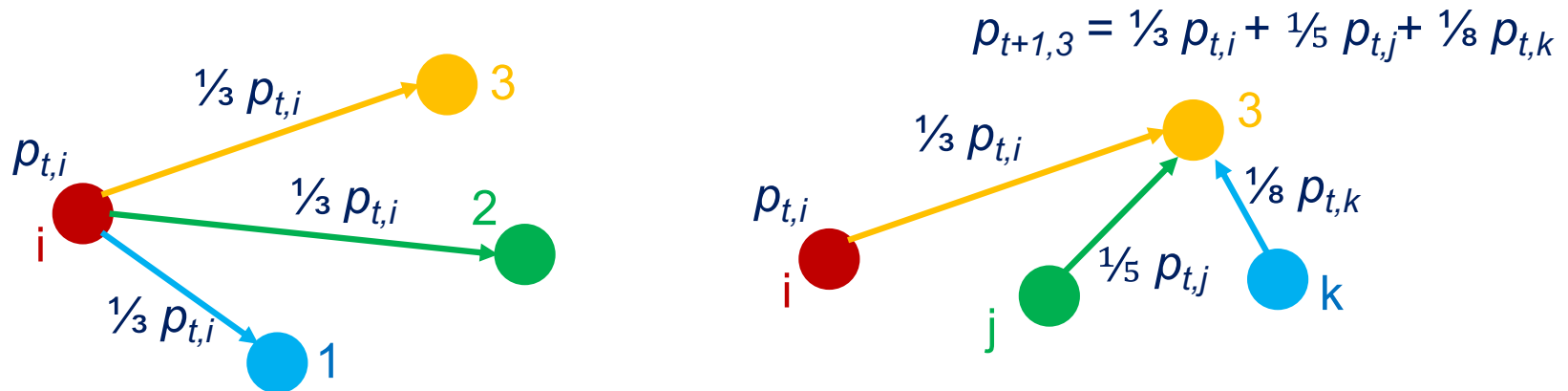




# A random walk on www

the rationale behind PageRank

- at time  $t$ , a web surfer is at page  $i$  with probability  $p_{t,i}$
- let the surfer choose with **equal probability** one of the sites linked by site  $i$

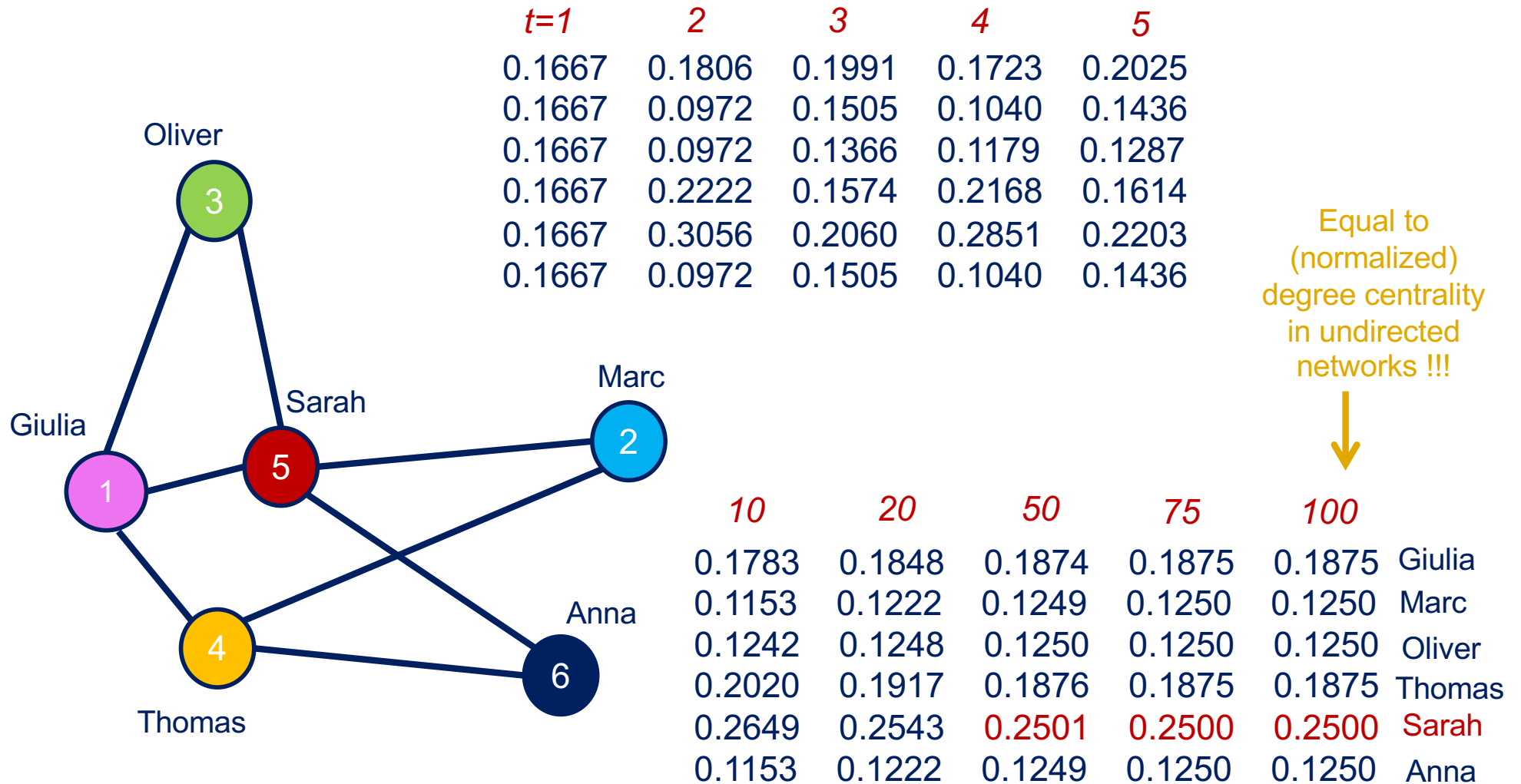


- this identifies a Markov chain
- after a while probabilities settle to a steady state = the PageRank vector



# Example

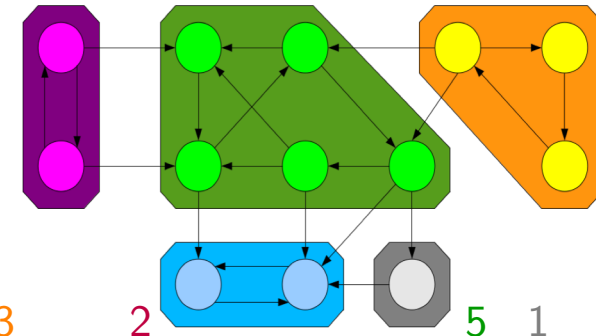
of the random walk effect on a friends' network





# Matrix formalization of the random walk

- $\mathbf{p}_{t+1} = \mathbf{M} \mathbf{p}_t$
- $\mathbf{p}_t$  stochastic vector  
(positive entries which sum up to 1)
- $\mathbf{M}$  **normalized** adjacency matrix (column stochastic)
- $\mathbf{M} = \mathbf{A} \text{diag}^{-1}(\mathbf{d})$
- $\mathbf{d} = \mathbf{A}^T \mathbf{1}$  output degree vector
- $\mathbf{p}_\infty = \mathbf{M} \mathbf{p}_\infty$  converges to an eigenvector of  $\mathbf{M}$  (with eigenvalue 1)
- $\mathbf{p}_\infty = \mathbf{d}$  for undirected networks where  $\mathbf{A} = \mathbf{A}^T$



$\mathbf{M} =$

	3	2	5			1	2	
	$\frac{1}{3}$							
	1							
		$\frac{1}{2}$						
	$\frac{1}{3}$	$\frac{1}{2}$				$\frac{1}{2}$		
	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$			
			$\frac{1}{2}$	1	$\frac{1}{3}$			
					$\frac{1}{3}$	$\frac{1}{3}$	0	
					$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{2}$	1
							1	1
							1	1

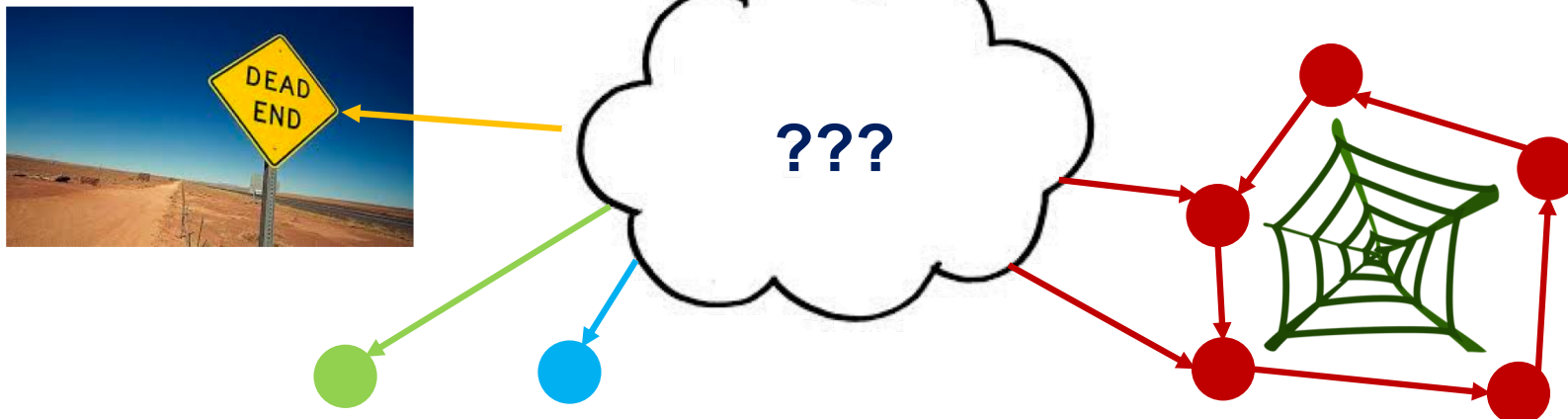
↑   ↑

**columns sum to 1**

62

With high probability the surfer ends in:

- ❑ **Dead ends**: some nodes do not have a way out  
= zero valued columns of  $M$
- ❑ **Spider traps**: some set of nodes do not have a way out, and further induce a **periodic** behaviour



## Idea:

❑ the surfer does not necessarily move to one of the links of the page she/he is viewing

❑ with a certain **probability**, might jump to a **random page**

❑  $p_{t+1} = c M p_t + (1-c) q$

↓  
damping factor, typically  $c = 0.85$ , meaning that **85% of the times** the surfer moves to one of the links of the page

↘  
the remaining  $1 - c = 15\%$  of the times the surfer moves at random according to a probability vector  $q$  independent of the node she/he is in, e.g.,  $q = 1/N$  for uniform probability





dead ends

$A_0$  ← original adjacency matrix  
(can be fractional)

no dead ends

$A = A_0 + b e^T$  ← teleportation  
← indicating vector of dead ends

normalization

$M = A \text{diag}^{-1}(d), \quad d = A^T \mathbf{1}$

no spider traps

$M_1 = c M + (1-c) q \mathbf{1}^T$   
equivalent formulation  
matrix is no more sparse

Markov chain

$p_{t+1} = M_1 p_t$

PageRank centrality vector

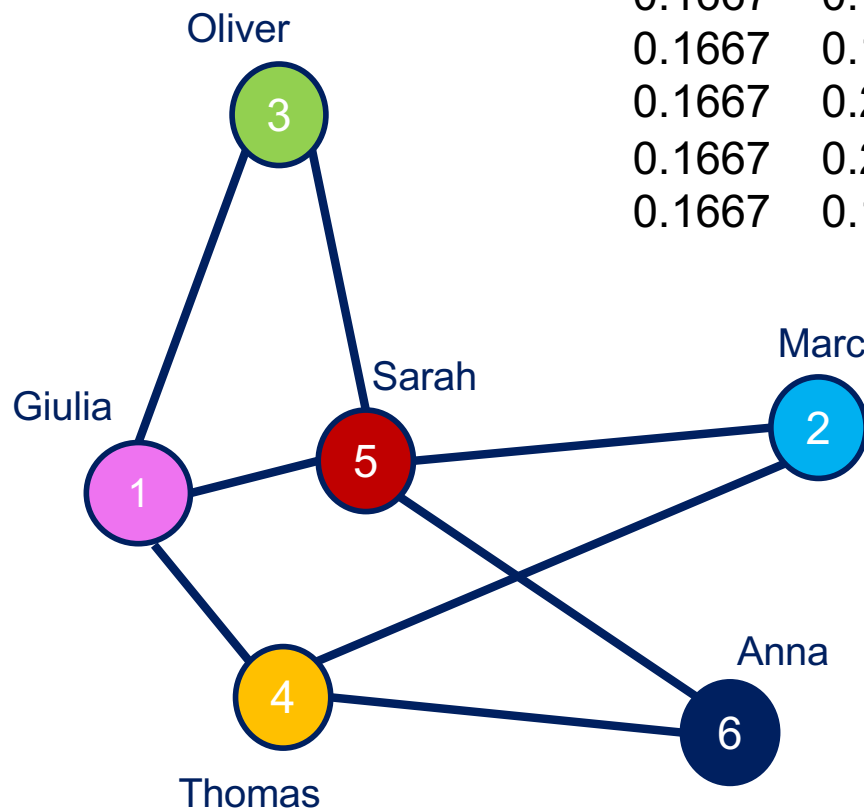
PageRank equation

$r = c M r + (1-c) q, \quad r = p_\infty$



# Example

of PageRank with restart on a friends' network



	<i>t=1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
Oliver	0.1667	0.1785	0.1919	0.1754	0.1912
Giulia	0.1667	0.1076	0.1461	0.1176	0.1382
Oliver	0.1667	0.1076	0.1361	0.1246	0.1302
Giulia	0.1667	0.2139	0.1671	0.2035	0.1746
Oliver	0.1667	0.2847	0.2128	0.2614	0.2276
Giulia	0.1667	0.1076	0.1461	0.1176	0.1382

not anymore  
identical to  
degree  
centrality !!!



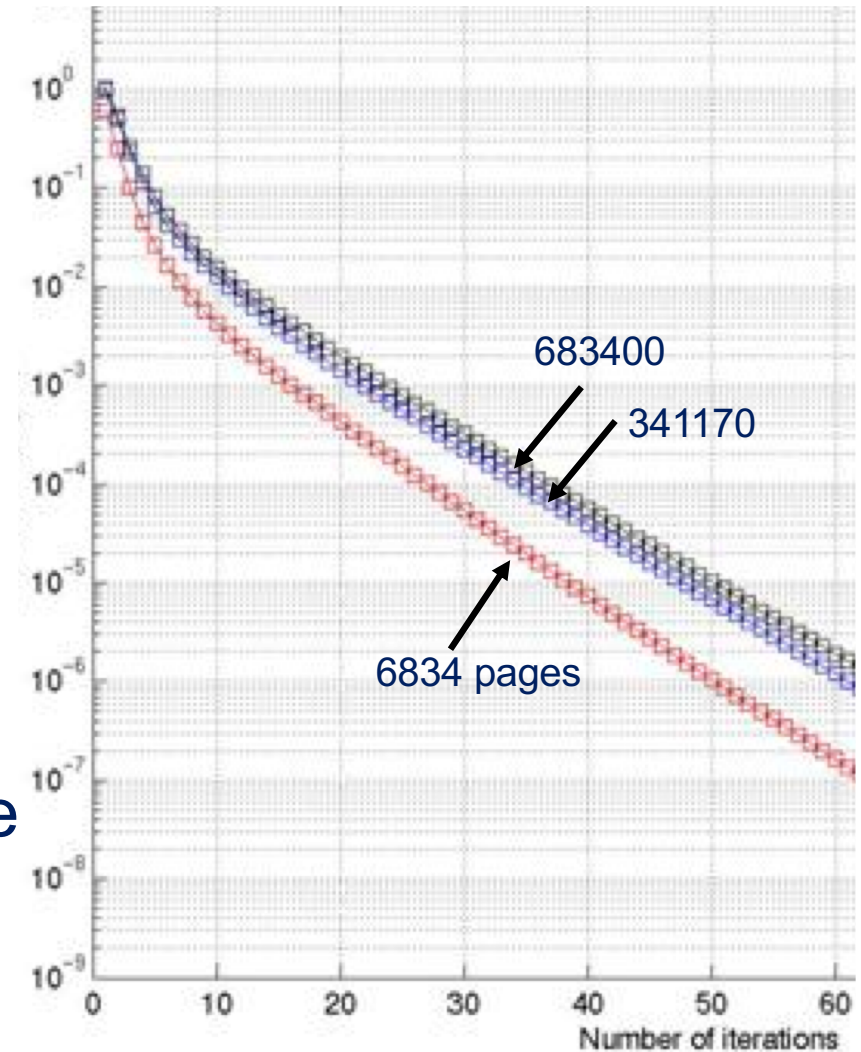
	<i>10</i>	<i>20</i>	<i>50</i>	<i>75</i>	<i>100</i>	
Giulia	0.1820	0.1839	0.1840	0.1840	0.1840	Giulia
Marc	0.1273	0.1293	0.1294	0.1294	0.1294	Marc
Oliver	0.1283	0.1285	0.1285	0.1285	0.1285	Oliver
Thomas	0.1902	0.1873	0.1871	0.1871	0.1871	Thomas
Sarah	0.2449	0.2419	0.2417	0.2417	0.2417	Sarah
Anna	0.1273	0.1293	0.1294	0.1294	0.1294	Anna



# Convergence properties of PageRank

an overview

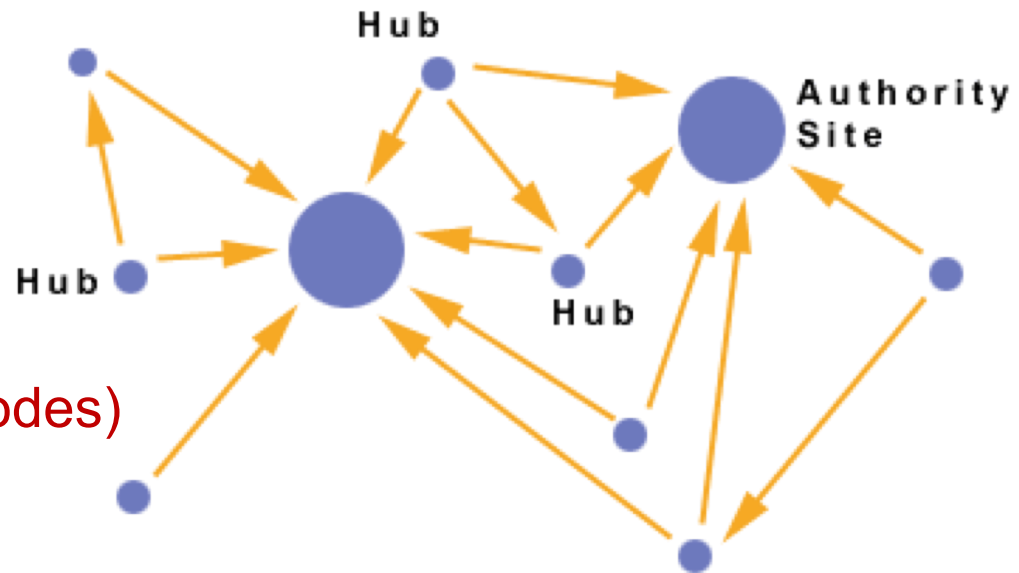
- The **PageRank** vector is the probability  $\mathbf{p}_t$  for large  $t$
- It corresponds to the **stationary behaviour** of the Markov chain
- $\mathbf{p}_\infty$  is **unique**
- $\mathbf{p}_\infty$  is a **stochastic vector** (with positive entries summing to 1)
- $\mathbf{p}_\infty$  depends on the choice of the teleportation vector  $\mathbf{q}$  (and of  $c$ )
- $\mathbf{p}_\infty$  converges in few iterations, typically  $\mathbf{p}_{40} \simeq \mathbf{p}_\infty$





## Authority (quality as a content provider)

nodes that contain useful information, or having a high number of edges pointing to them (e.g., course homepages)  
= PageRank vector  
(related to the in-degree of nodes)



## Hub (quality as an expert)

trustworthy nodes, or nodes that link to many authorities (e.g., course bulletin)  
= PageRank vector starting from  $\mathbf{A}_0^T$   
(related to the out-degree of nodes)

authority or hub?



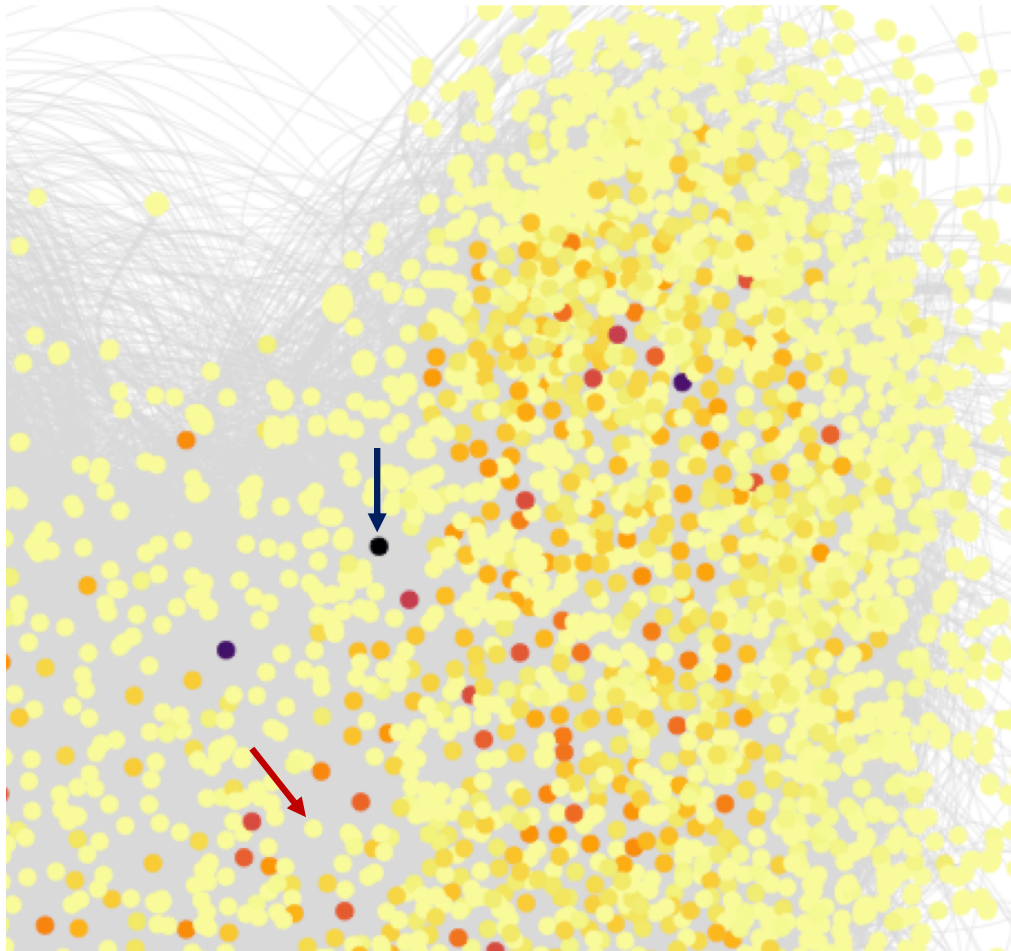


# Example of PageRank centrality

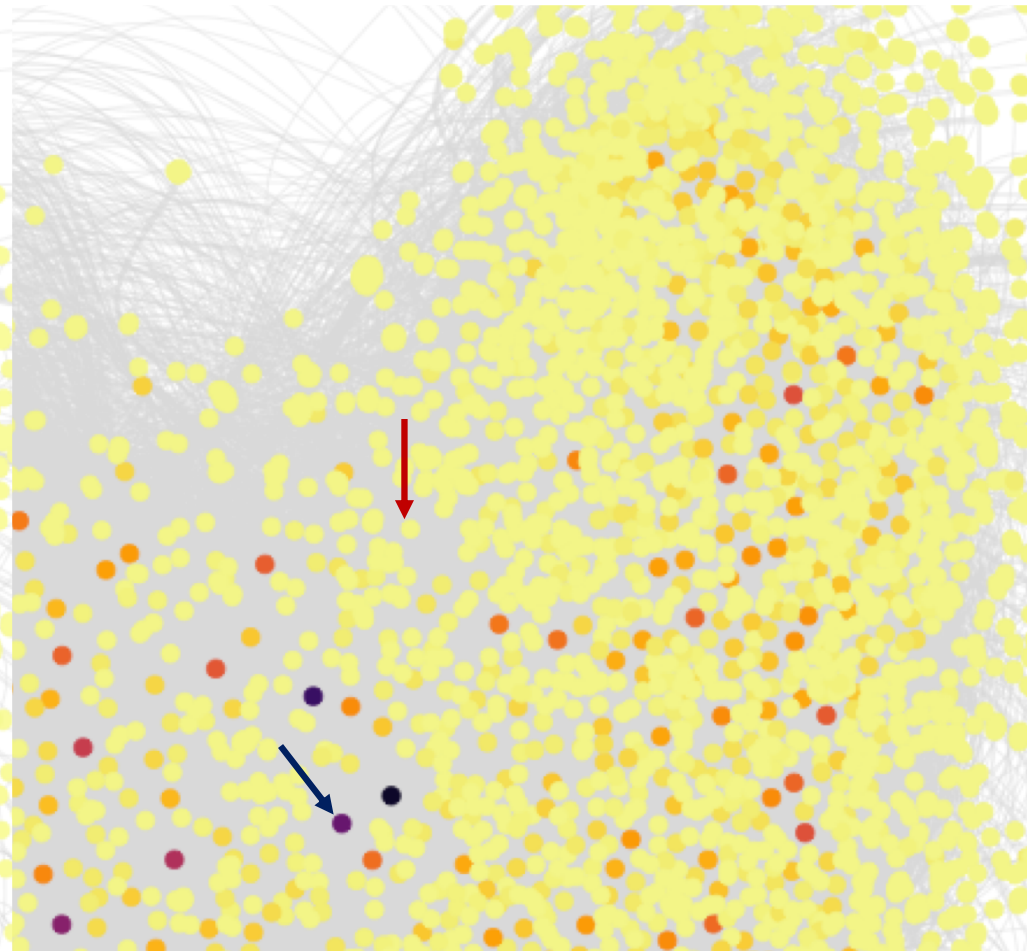
wikipedia administrator elections and vote history data

<https://snap.stanford.edu/data/wiki-Vote.html>

## Authorities

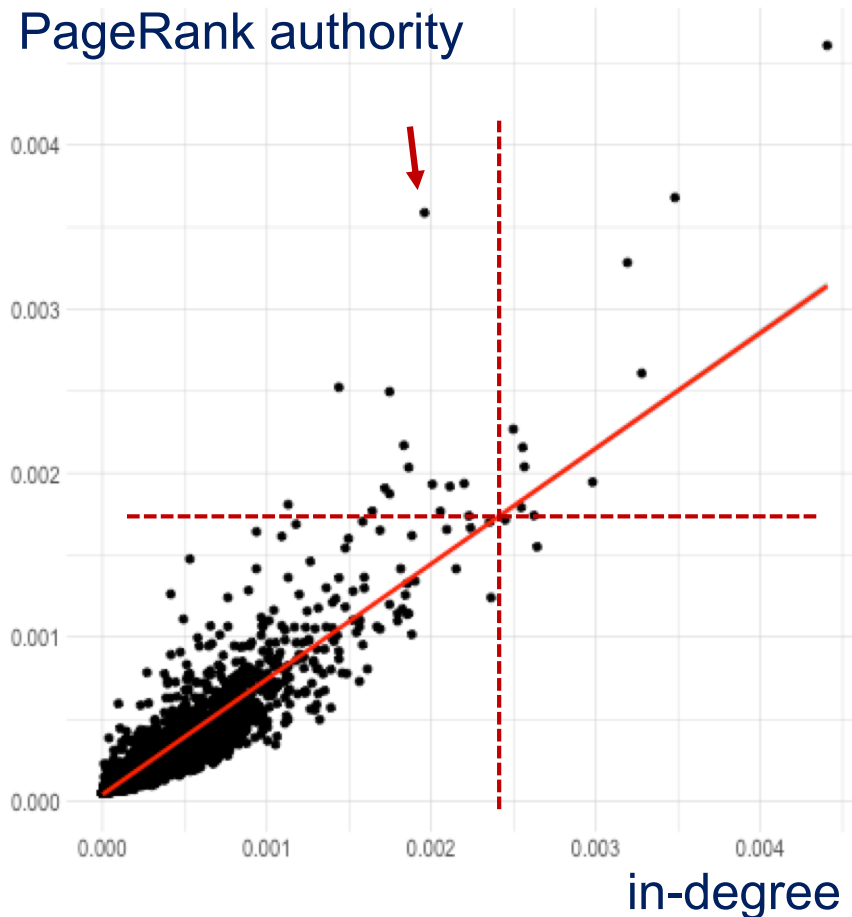


## Hubs

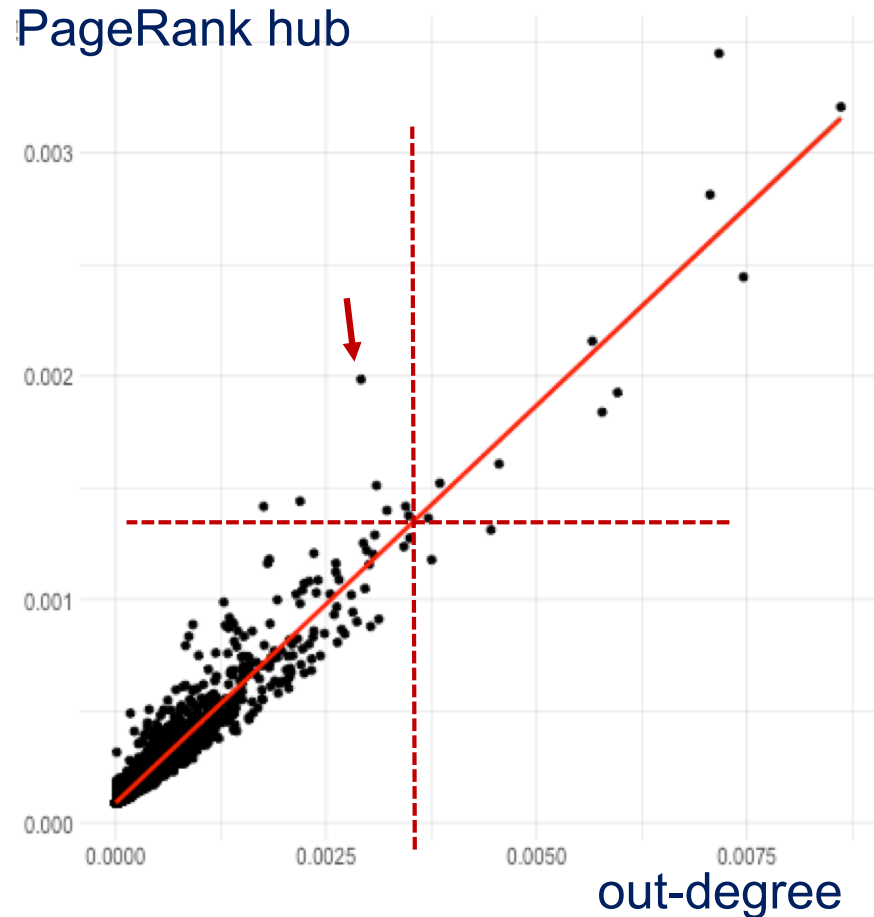




## Authorities



## Hubs

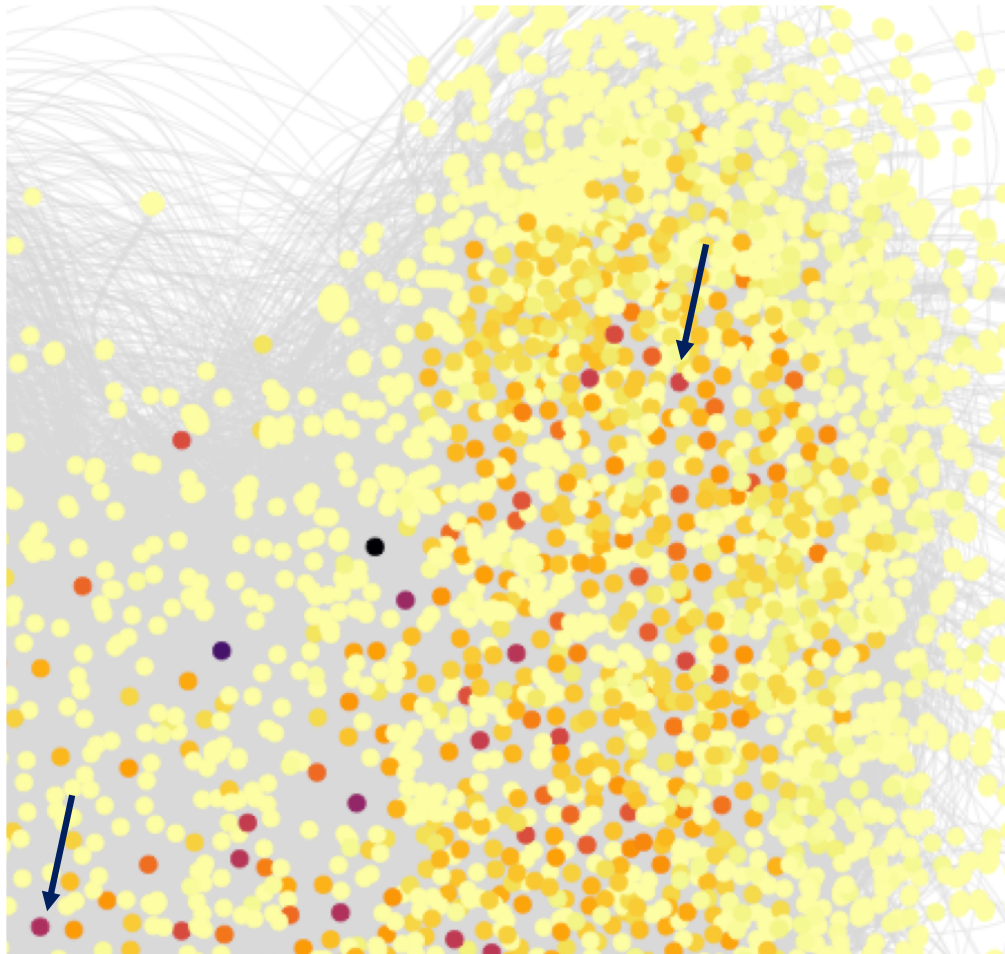




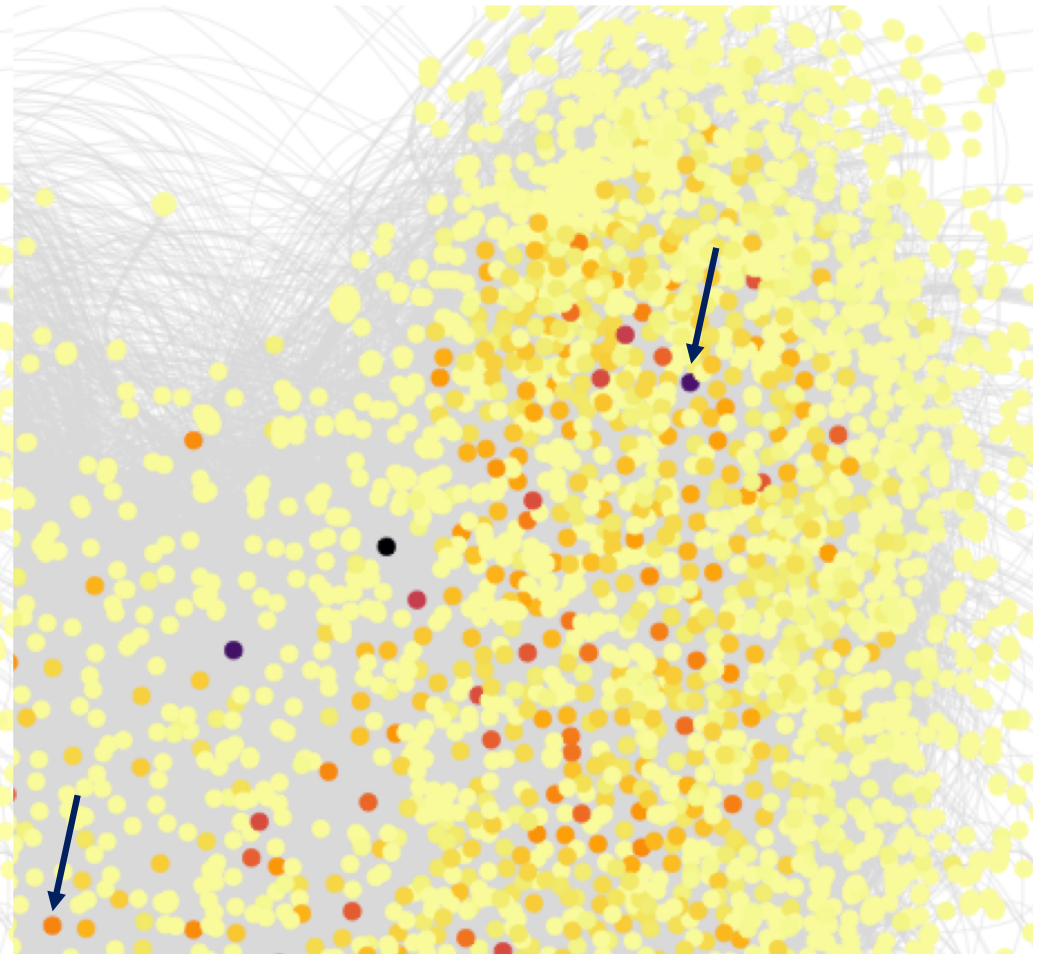
# PageRank versus degree authorities

wikipedia administrator elections and vote history data

## Degree



## PageRank

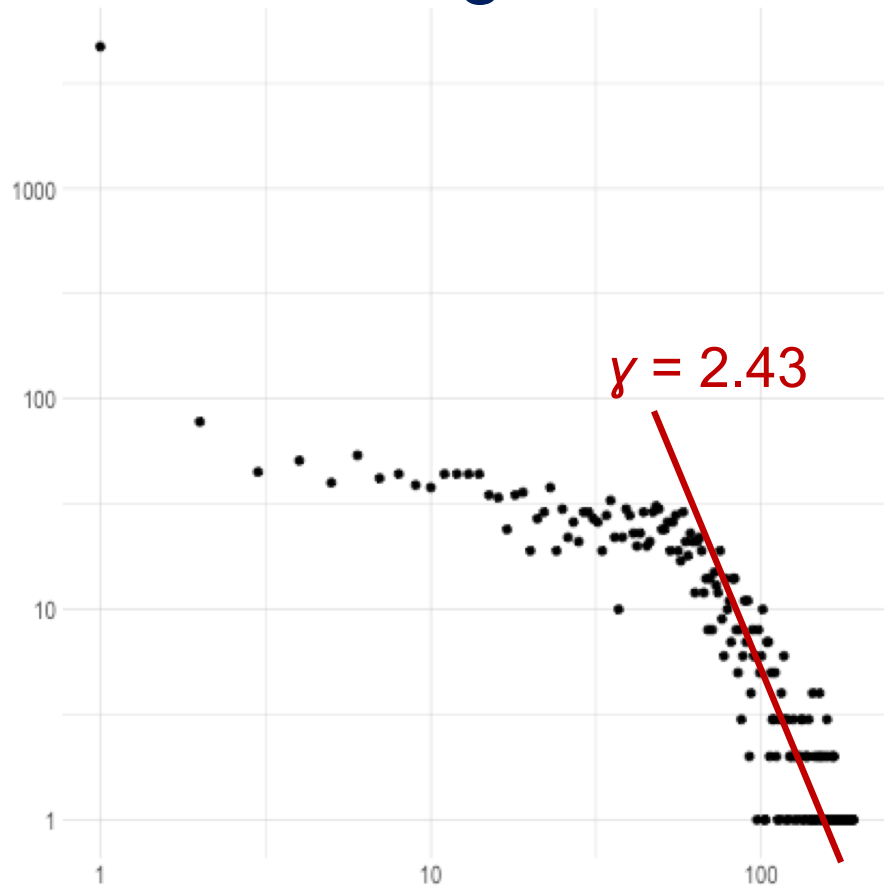




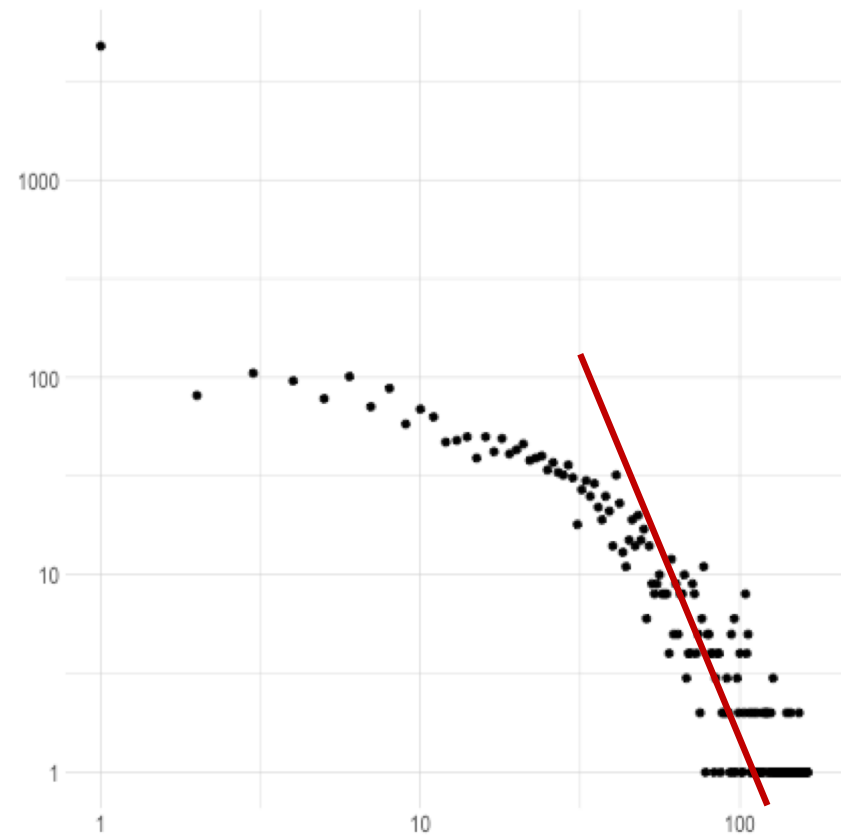
# PageRank versus degree authorities

wikipedia administrator elections and vote history data

## Degree



## PageRank

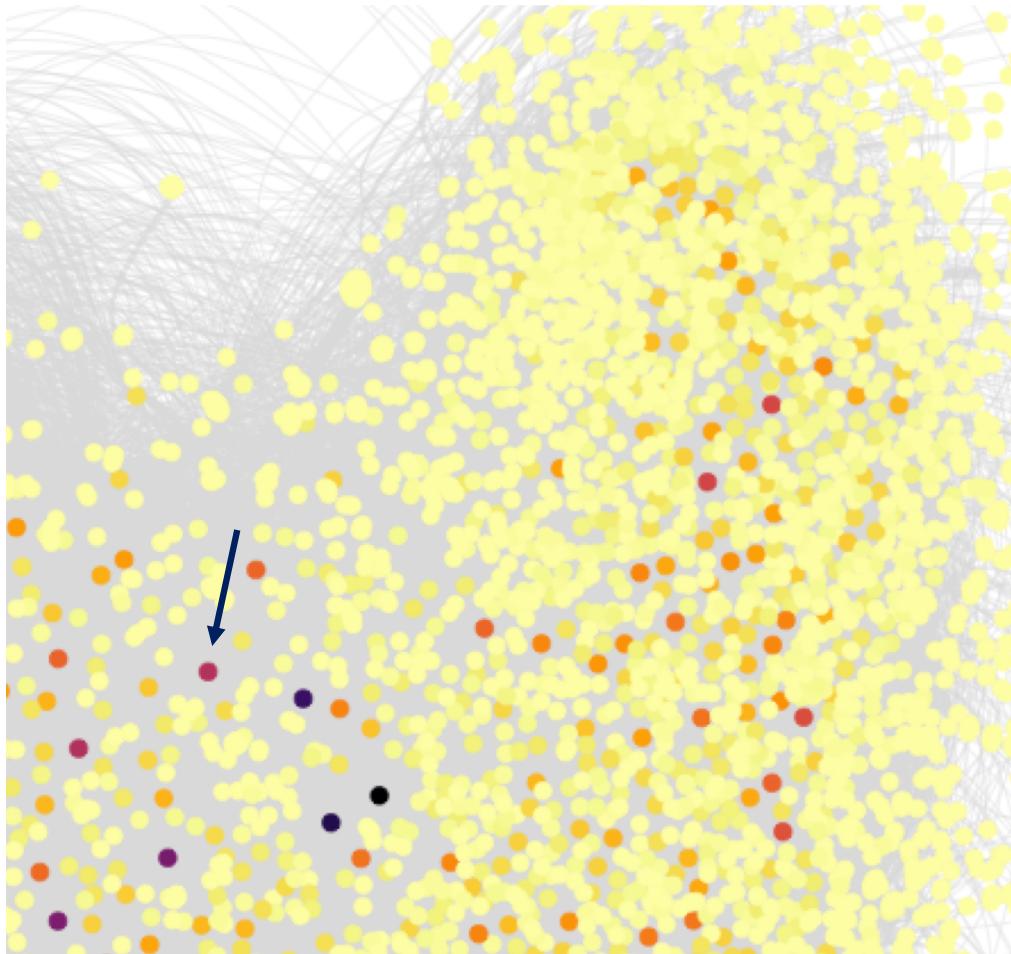




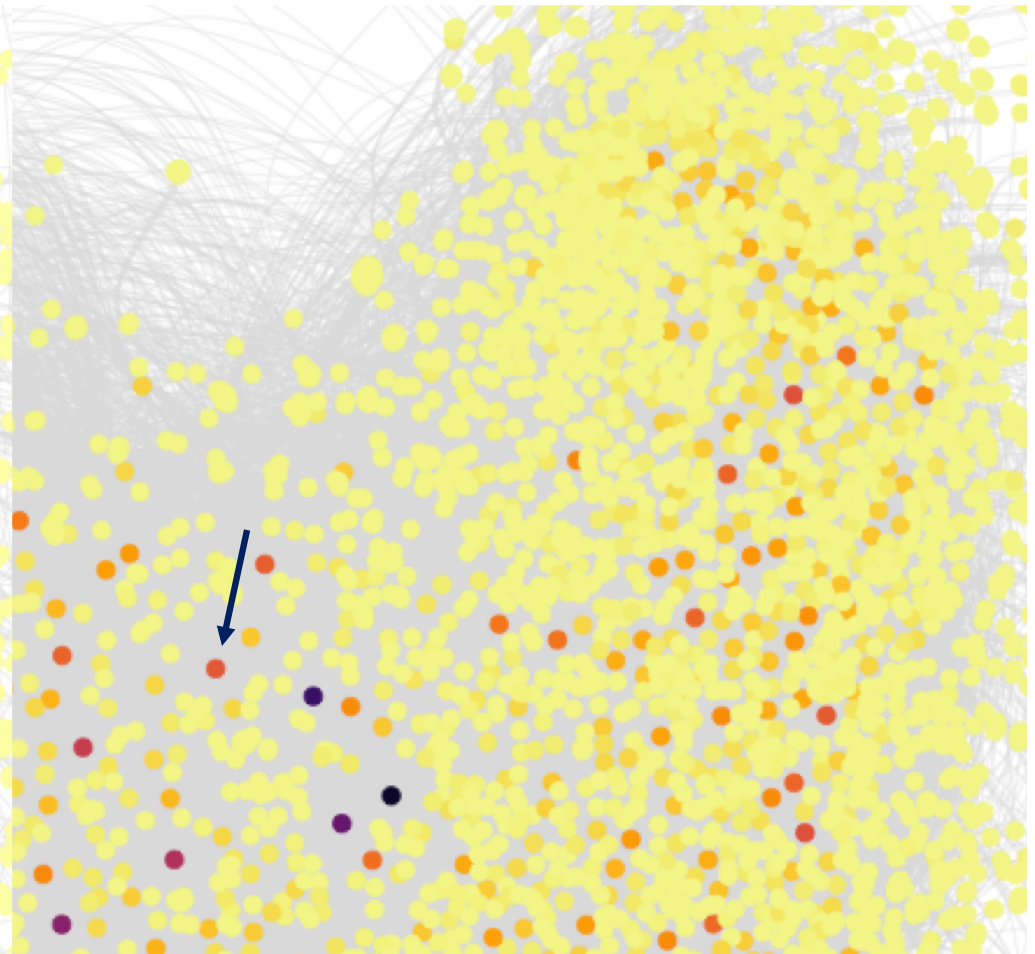
# PageRank versus degree hubs

wikipedia administrator elections and vote history data

## Degree

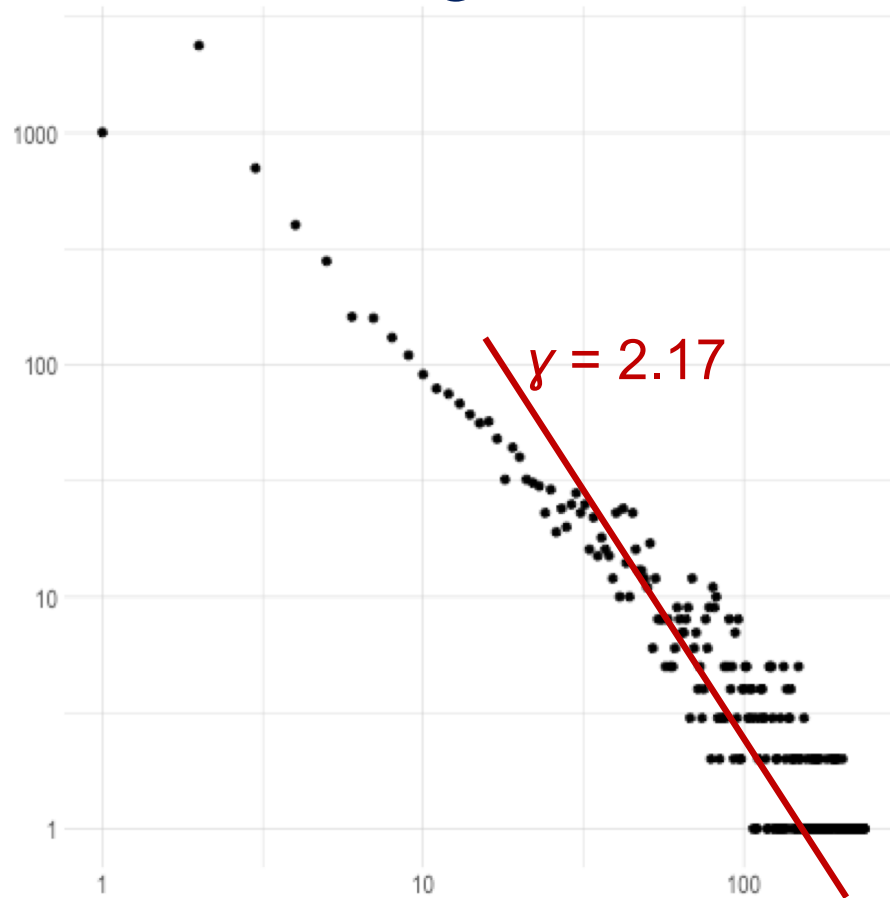


## PageRank

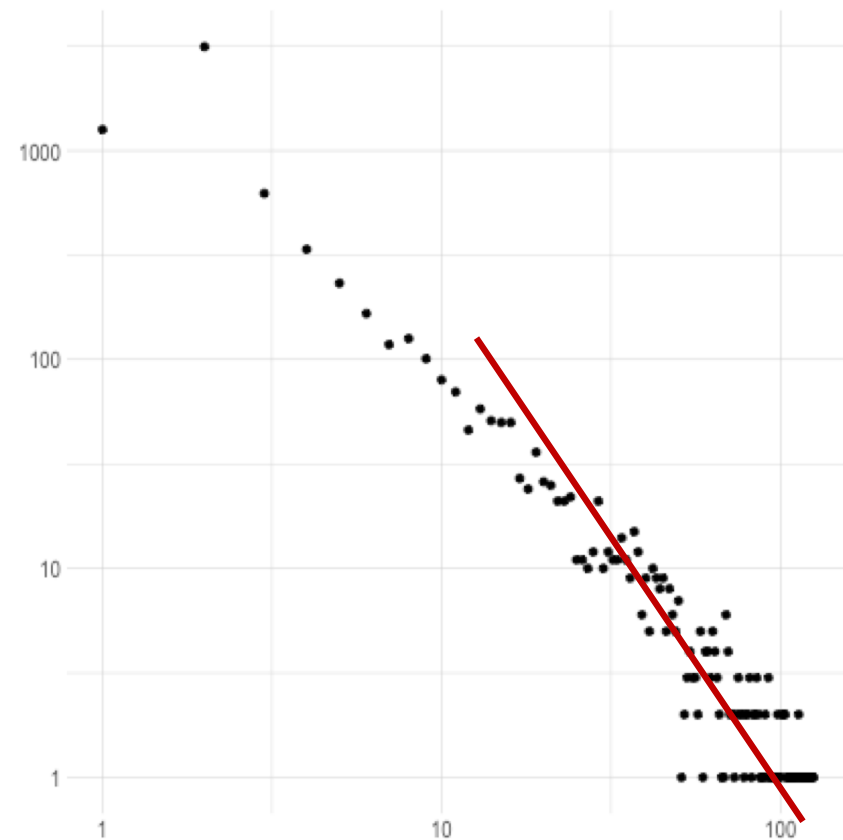




## Degree



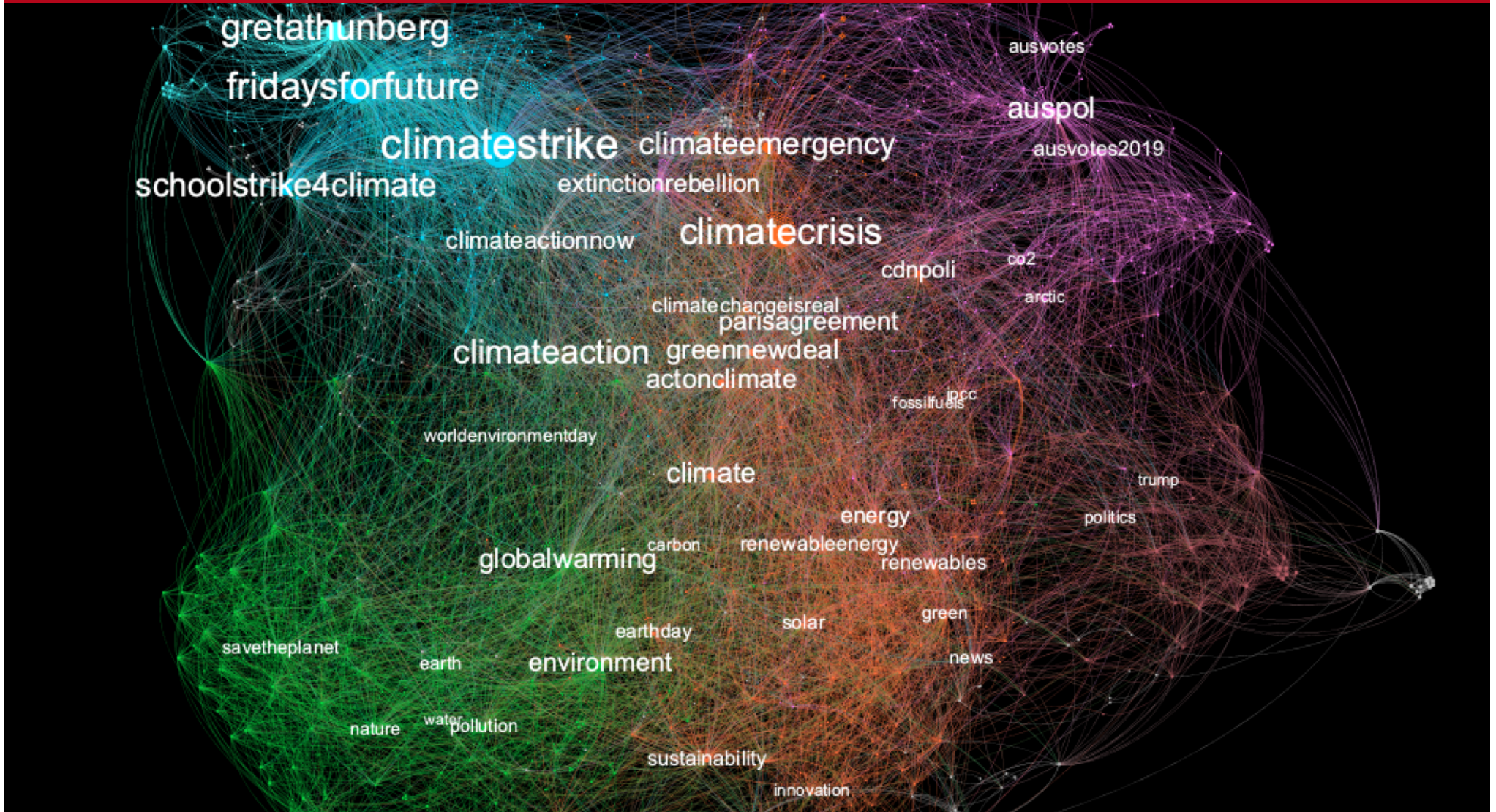
## PageRank





# PageRank on a semantic network

2019 hashtag network related to #climatechange  
(from Twitter, after #gretathunberg)







- ❑ Brin and Page, “The anatomy of a large-scale hypertextual web search engine,” 1998
- ❑ Page, Brin, Motwani, Winograd, “The PageRank Citation Ranking: Bringing Order to the Web,” 1999

<http://ilpubs.stanford.edu/422/1/1999-66.pdf>

  
<https://scholar.google.com/>

# Convergence properties

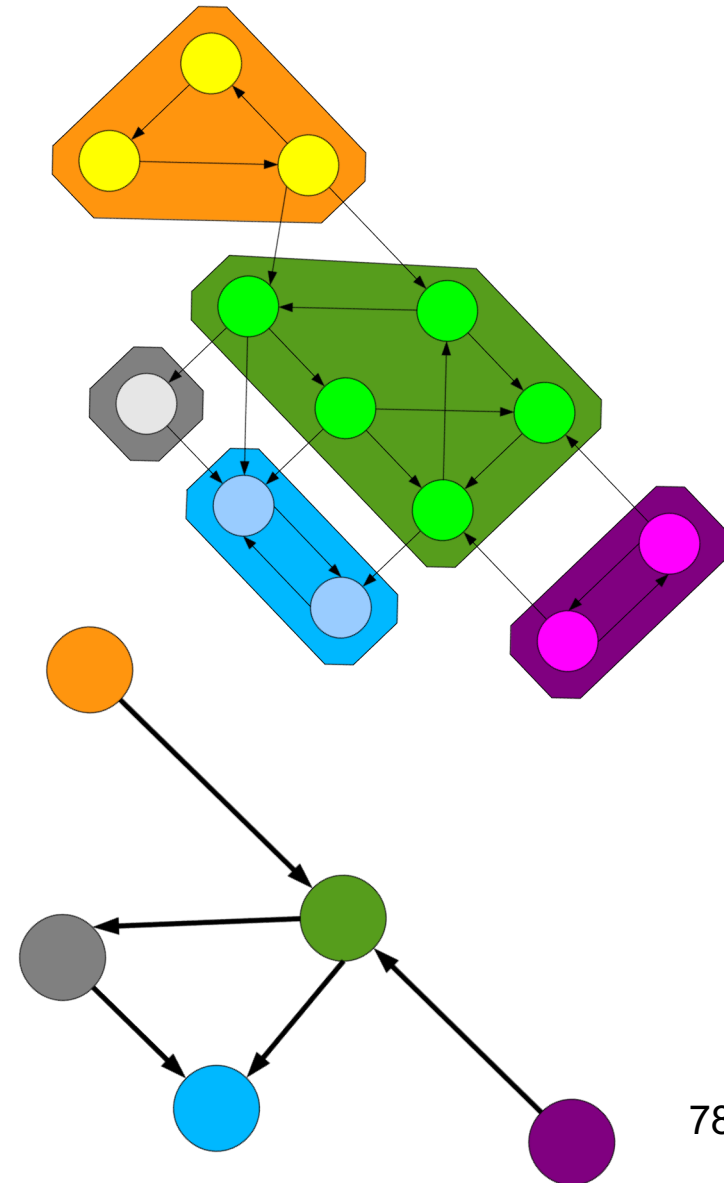
of PageRank power iterations



# The condensation graph

ordering an adjacency matrix

- Strong connectivity induces a **partition** in disjoint **strongly connected** sets  $\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_K$
- By reinterpreting the sets as nodes we obtain a **condensation graph**  $\mathcal{G}^*$  where  $i \rightarrow j$  is an edge if a connection exists between sets  $\mathcal{V}_i \rightarrow \mathcal{V}_j$





# Properties of the condensation graph

ordering an adjacency matrix

- $\mathcal{G}^*$  does not contain **cycles**

otherwise the sets in the cycle would be strongly connected

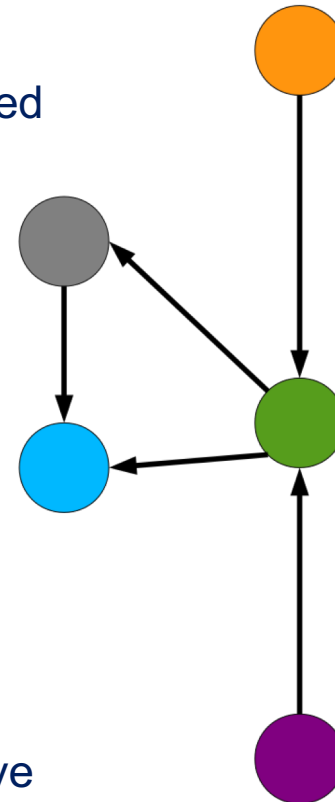
- $\mathcal{G}^*$  has at least one **root** and one **leaf**

and every node in the graph can be reached from one of the roots

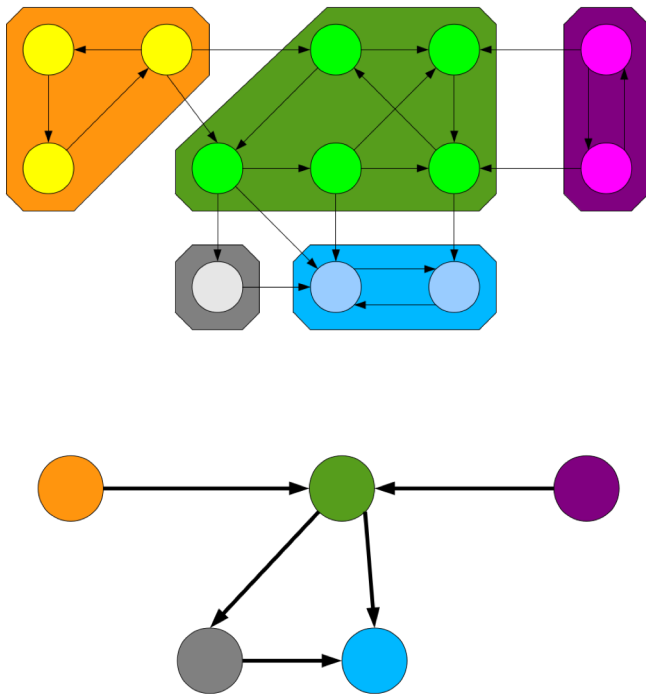
- $\mathcal{G}^*$  allows a particular **reordering**

where node  $n_i$  does not reach any of the nodes  $n_j$  with  $j < i$

procedure: identify a root  $n_1$  and remove it from the network, then identify a new root; cycle until all nodes have been selected



The **condensation graph** ordering induces a block-lower-triangular matrix structure on the adjacency matrix



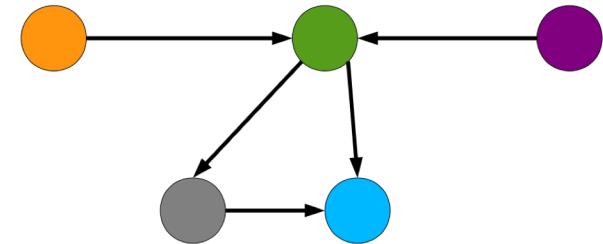
$M =$

	3	2	5	1	2
$\frac{1}{3}$	1				
1		$\frac{1}{2}$			
$\frac{1}{3}$		$\frac{1}{2}$		$\frac{1}{2}$	
$\frac{1}{3}$		$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	
		$\frac{1}{2}$	1	$\frac{1}{3}$	
			$\frac{1}{3}$	$\frac{1}{3}$	0
			$\frac{1}{3}$	$\frac{1}{3}$	1
				$\frac{1}{2}$	1

blocks in the diagonal are **irreducible** = no block-diagonal form ! 80

# Perron-Frobenius theorem of the condensation graph

the eigenvalues of the diagonal blocks, except for the leaves, lie inside the unit circle, i.e.,  $|\lambda| < 1$



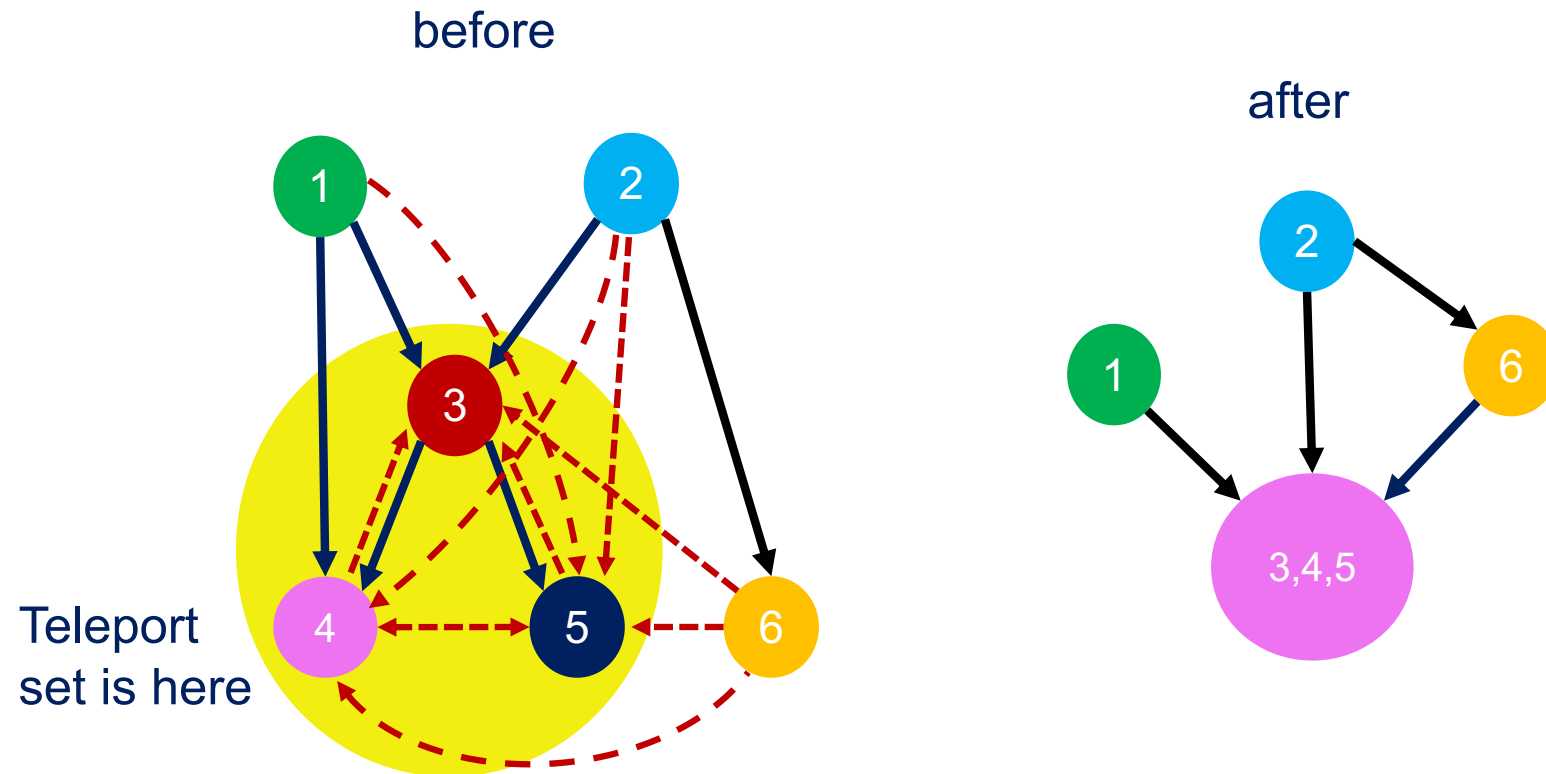
$M =$

	3	2	5	1	2
1	$\frac{1}{3}$ 1				
		$\frac{1}{2}$ $\frac{1}{2}$			
$\frac{1}{3}$ $\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$ $\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{2}$	
		$\frac{1}{2}$	1 $\frac{1}{3}$	$\frac{1}{3}$	
			$\frac{1}{3}$ $\frac{1}{3}$	0	
			$\frac{1}{3}$ $\frac{1}{3}$	1	1
			$\frac{1}{2}$	1	1

each leaf-block has at least one eigenvalue in the unit circle;  $\lambda=1$  is always available, the others are distinct

# The teleportation effect

it implies only one leaf



Hence  $M_1$  carries only **one** eigenvector associated with the eigenvalue  $\lambda=1$



- ❑ PageRank matrix  $M_1 = c M + (1-c) q \mathbf{1}^T$
- ❑ Normalization property  $\mathbf{1}^T M_1 = \mathbf{1}^T$
- ❑ Jordan form  $M_1 = V J V^{-1}$

carries the right (generalized)  
eigenvectors  $\mathbf{e}_i$  of  $M_1$

carries the  
eigenvalues of  $M_1$

$$J = \begin{bmatrix} \boxed{\begin{matrix} \lambda_1 & 1 \\ & \lambda_1 & 1 \\ & & \lambda_1 \end{matrix}} & & & \\ & \boxed{\begin{matrix} \lambda_2 & 1 \\ & \lambda_2 \end{matrix}} & & \\ & & \boxed{\lambda_3} & \\ & & & \ddots \\ & & & & \boxed{\begin{matrix} \lambda_n & 1 \\ & \lambda_n \end{matrix}} \end{bmatrix}$$

$$\begin{aligned} \mathbf{1}^T M_1 V &= \mathbf{1}^T V \\ &= \mathbf{1}^T V J \end{aligned} \quad \Rightarrow \quad \underbrace{\mathbf{1}^T V}_{\rho} \underbrace{(J - I)}_{\text{only one value is 0}} = 0$$

Hence  $\mathbf{1}^T \mathbf{e}_i = 0$  for  $i > 1$ , i.e., except for the eigenvector associated with eigenvalue 1





$$\underbrace{M_1 e_i = c M e_i}_{\text{same eigenvalues of } M, \text{ but multiplied by } c !!!} + (1-c) \cancel{q 1^T e_i} \quad \text{for } i > 1$$

same eigenvalues of  $M$ ,  
but multiplied by  $c$  !!!



- $M_1$  has **one** eigenvalue equal to **1**
- The remaining eigenvalues satisfy  $|\lambda| \leq c$

Haveliwala and Kamvar, "The second eigenvalue of the Google matrix," 2003

<http://ilpubs.stanford.edu:8090/582/1/2003-20.pdf>



$$\mathbf{p}_t = \mathbf{M}_1 \mathbf{p}_{t-1} = \mathbf{M}_1^t \mathbf{p}_0 = \mathbf{V} \mathbf{J}^t \mathbf{V}^{-1} \mathbf{p}_0$$

gets large for high multiplicity

max eigenvalue multiplicity

$$\square \|\mathbf{p}_t - \mathbf{p}_\infty\|_2 \lesssim K c^t t^{m-1} \sim K c^t$$

$$\square \text{Triangular inequality: } \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2 \lesssim 2K c^t$$

$$\square \text{Precision } \varepsilon: \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2 < \varepsilon$$

$$\square \text{Iterations required: } t = \lceil [\ln(2/\varepsilon) + \ln(K)] / \ln(1/c) \rceil$$

precision  $10^{-3} \rightarrow 7.6$

$c=0.85 \rightarrow 1/\ln(1/c) = 6$

Is usually small  
 $\rightarrow$  fast algorithm

# Local PageRank

measuring similarity/closeness among nodes



# Measuring closeness: LocalPageRank

for the eigenstructure of the PageRank matrix

## Idea

- ❑ Measure **similarity** / closeness to node  $i$  by applying PageRank with teleport set  $S=\{i\}$ , i.e., with  $\mathbf{q} = \delta_i$

## Result

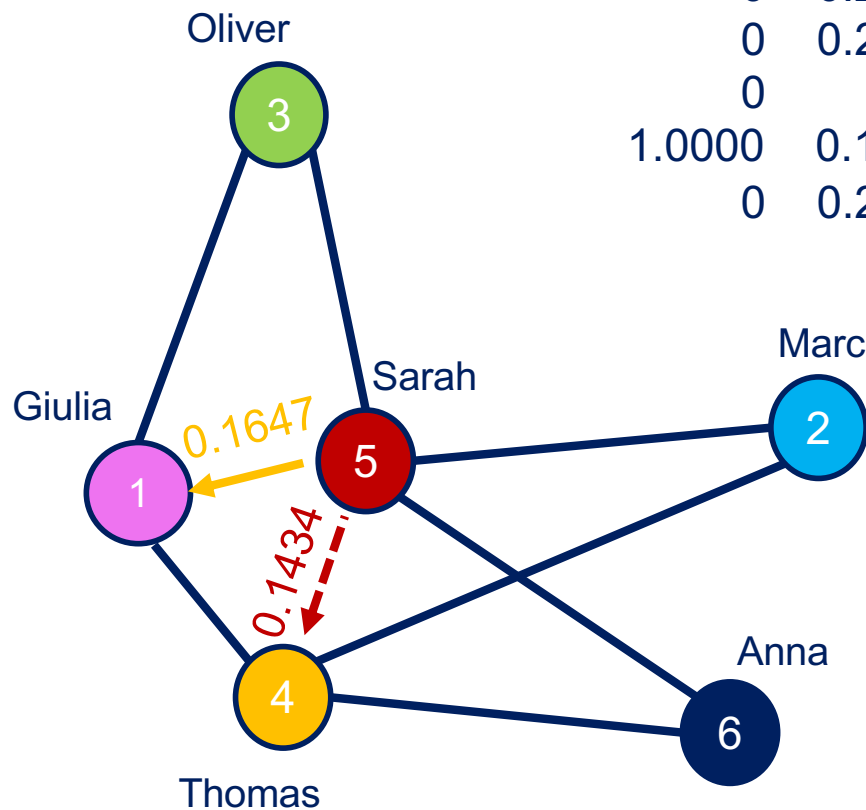
- ❑ Measures direct and indirect multiple connections, their quality, degree or weight





# Example

who's Sara's best friend?



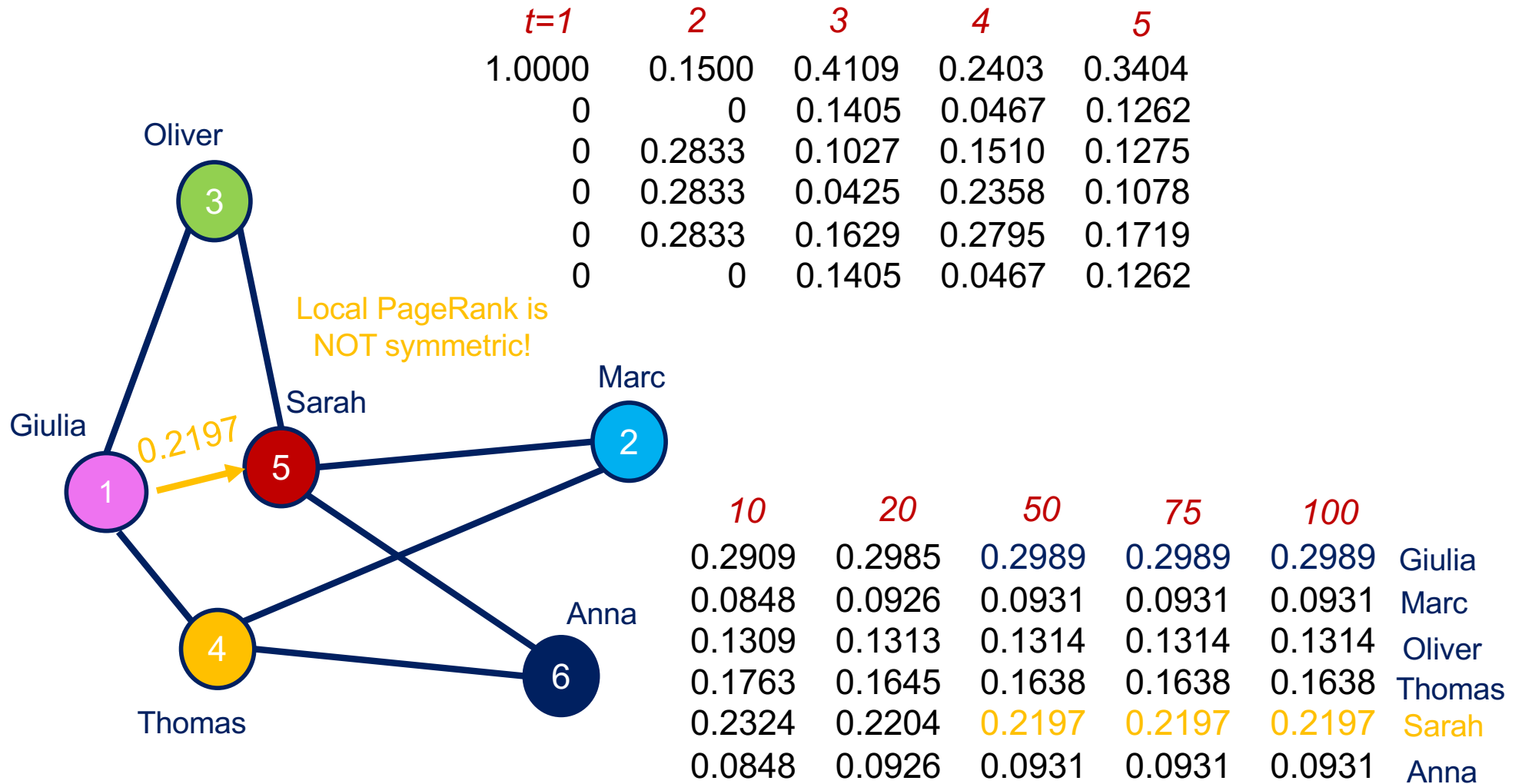
	<i>t=1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
0	0	0.2125	0.1222	0.2096	0.1290
0	0	0.2125	0.0319	0.1705	0.0708
0	0	0.2125	0.0921	0.1369	0.1127
0	0	0	0.2408	0.0617	0.2043
1.0000	0.1500	0.4811	0.2508	0.4125	
0	0.2125	0.0319	0.1705	0.0708	

	<i>10</i>	<i>20</i>	<i>50</i>	<i>75</i>	<i>100</i>	
0.1743	0.1653	0.1647	0.1647	0.1647	0.1647	Giulia
0.1238	0.1144	0.1138	0.1138	0.1138	0.1138	Marc
0.1206	0.1199	0.1199	0.1199	0.1199	0.1199	Oliver
0.1285	0.1426	0.1434	0.1434	0.1434	0.1434	Thomas
0.3290	0.3435	0.3444	0.3444	0.3444	0.3444	Sarah
0.1238	0.1144	0.1138	0.1138	0.1138	0.1138	Anna



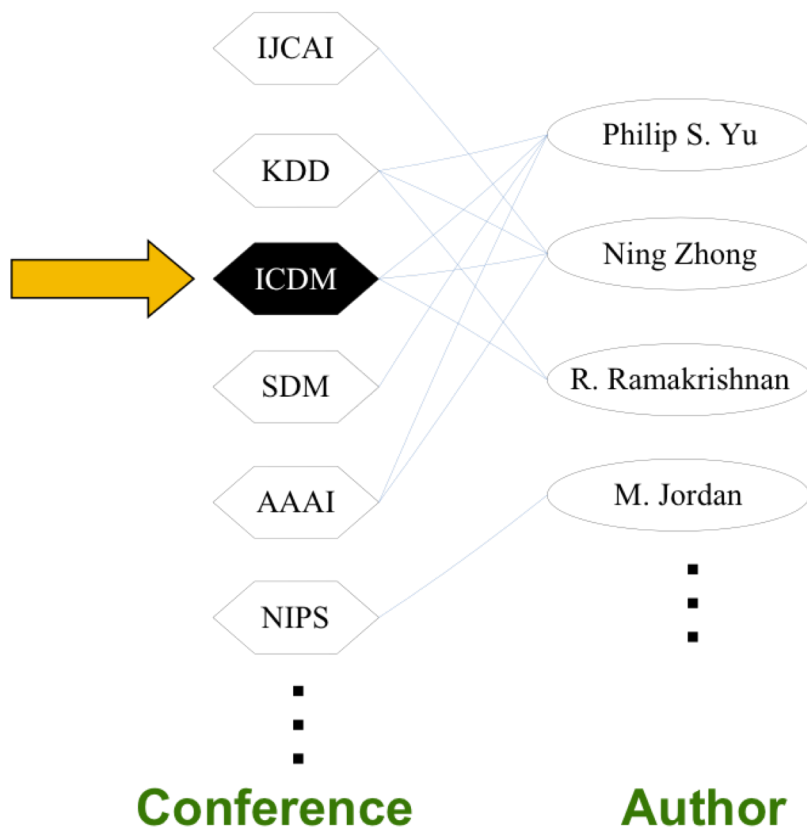
# Example

who's Giulia's best friend?

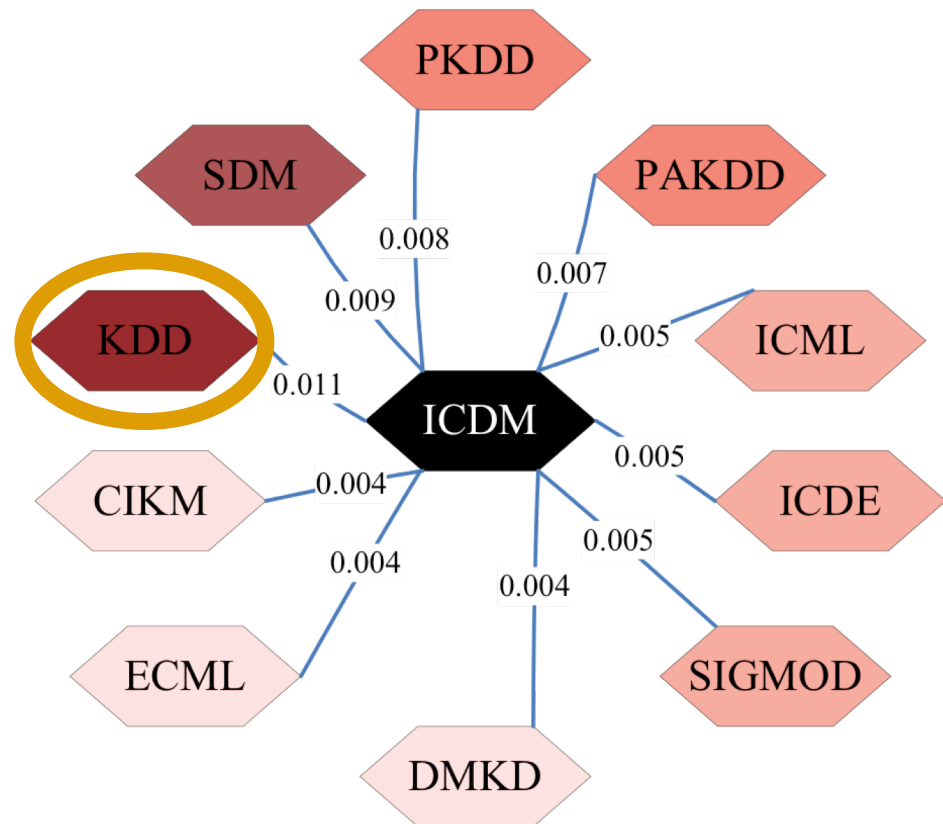




what is the most related conference to ICDM?



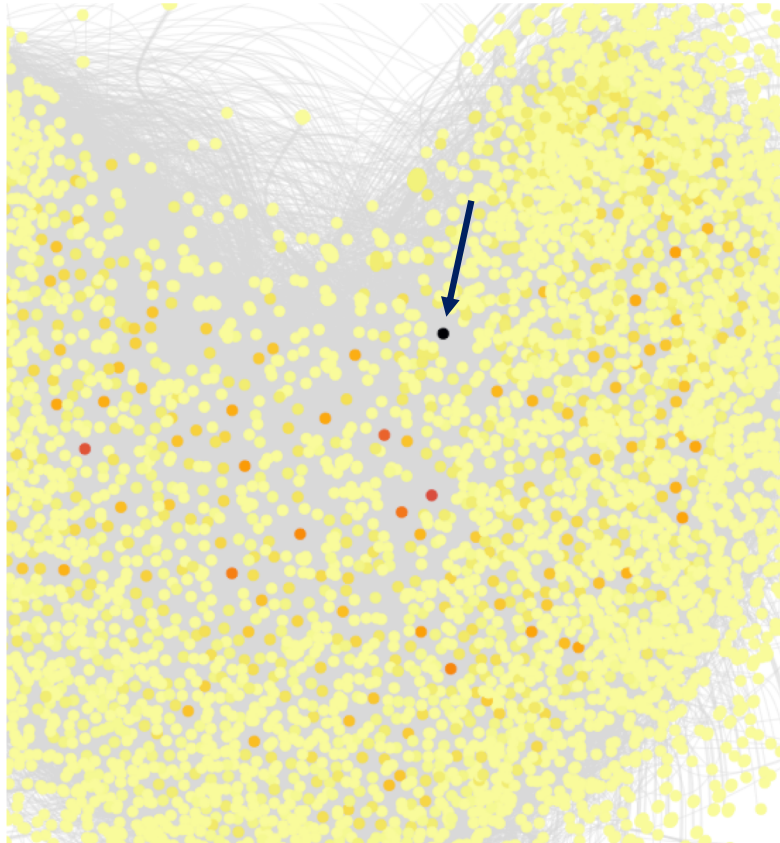
## Top 10 ranking results



ICDM = international conf. on data mining  
KDD = knowledge discovery and data mining

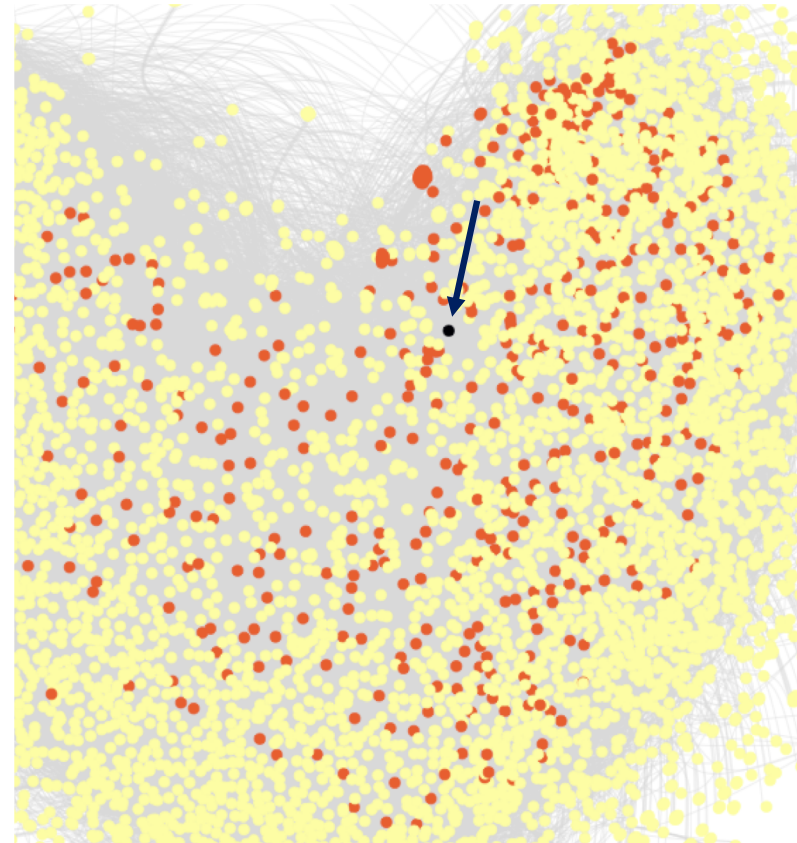


## Local PageRank



neighbours **authority score** =  
local node  $\rightarrow$  neighbours

## 1-hop out-neighbours







UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# On the complexity of Local PageRank

approximate PageRank

Andersen, Chung, Lang, “Local graph partitioning using PageRank vectors,” 2006

<https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=4031383>

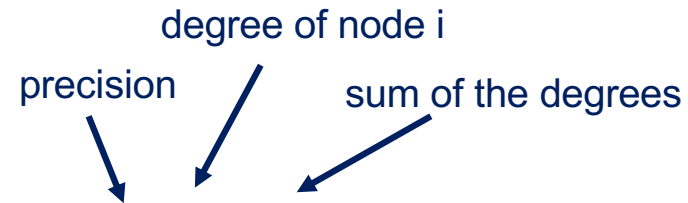
use [institutional Sign In](#) with your unipd credentials



# Approximate PageRank algorithm

the push operation

- Start from  $\mathbf{u} = \mathbf{0}$  and  $\mathbf{v} = \mathbf{q}$
- To all the nodes  $i$  satisfying  $\underline{v_i} > \varepsilon \underline{d_i/D}$  apply the **push** operation



$\mathbf{u}^+ = \mathbf{u} + (1-c) \delta$

$\mathbf{v}^+ = \mathbf{v} - \delta + c \mathbf{M} \delta$

u constantly increases

v always positive

only one active element in position i with value  $v_i$

Returns  $\mathbf{u} \simeq \mathbf{r}$  with precision  $|\mathbf{r} - \mathbf{u}|_1 < \varepsilon$   
It is simple



column stochastic matrix  $\mathbf{1}^T \mathbf{M} = \mathbf{1}^T$

□ PageRank equation  $\mathbf{r}_q = c \mathbf{M} \mathbf{r}_q + (1-c) \mathbf{q}$

stochastic ranking vector  
 $\mathbf{1}^T \mathbf{r}_q = 1, \mathbf{r}_q \geq 0$

stochastic Teleport vector  
 $\mathbf{1}^T \mathbf{q} = 1$

□ Alternative equation  $\mathbf{r}_q = (\mathbf{I} - c \mathbf{M})^{-1} (1-c) \mathbf{q}$

linear in  $\mathbf{q}$

$$\mathbf{r}_{a\mathbf{u}+b\mathbf{v}} = a \mathbf{r}_u + b \mathbf{r}_v$$



# Modifying the PageRank equation

the lemma for the proof

one-step random walk

□ PageRank equation  $r_q = c r_{Mq} + (1-c) q$



□  $r_q = (I - c M)^{-1} (1-c) q$

□  $r_q = (1-c) \sum (c M)^k q$

□  $M r_q = (1-c) \sum (c M)^k M q$

□  $M r_q = r_{Mq}$



# Main property of push: $r_q = u + r_v$ almost there

- At starting point  $u = 0$  and  $v = q$  imply  $r_q = 0 + r_q$
- The following steps are proved by induction

$$u^+ = u + (1-c) \delta$$
$$v^+ = v - \delta + c M \delta$$



$$u^+ + r_{v^+} = u + (1-c) \delta + \overbrace{r_v - r_\delta + c r_{M\delta}}^{\text{by linearity}}$$
$$r_\delta - (1-c) \delta$$



$$u^+ + r_{v^+} = u + r_v = r_q$$



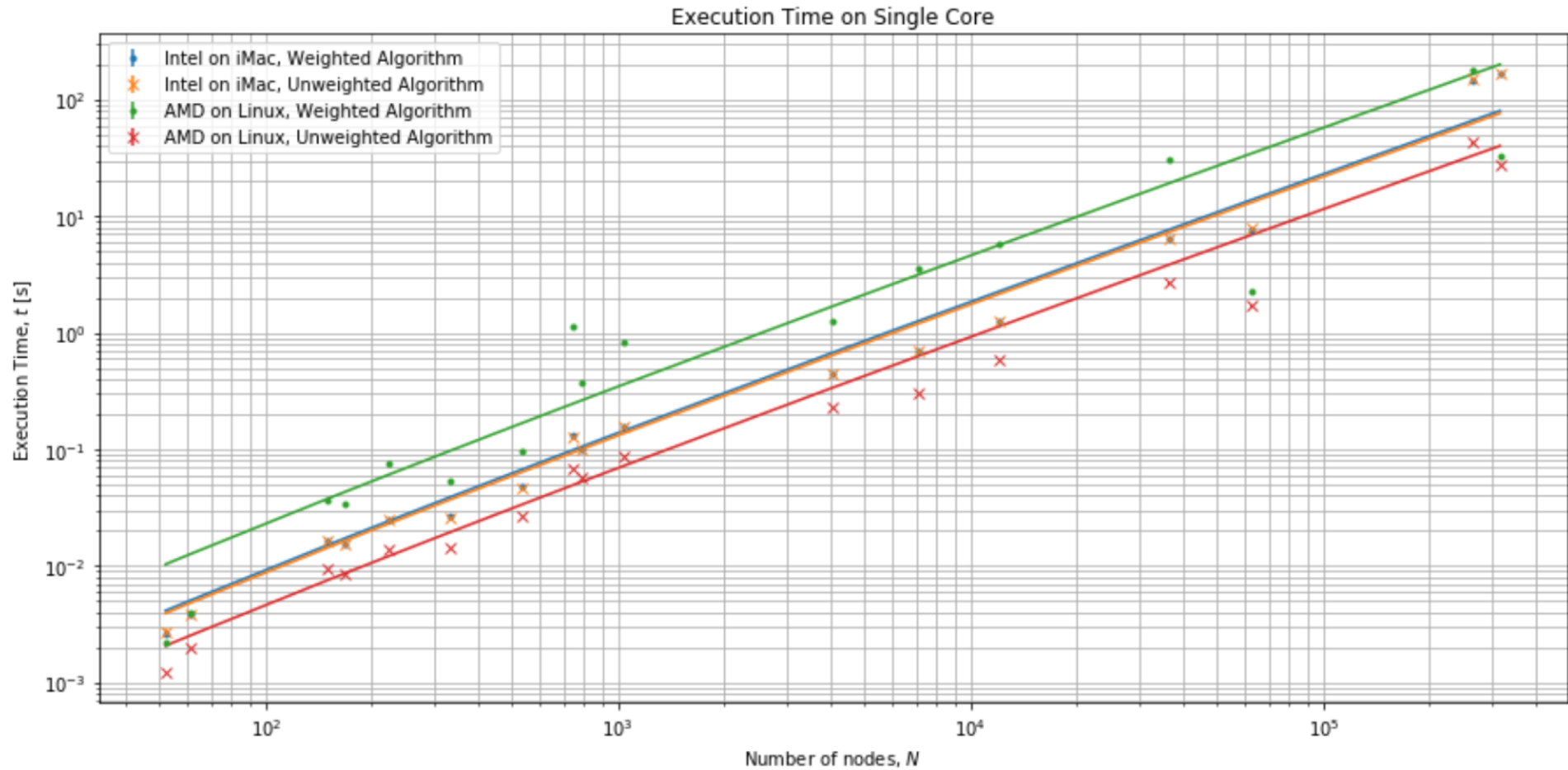
- The push property implies  $\mathbf{r}_q = \mathbf{u} + \mathbf{r}_v$
- Hence  $\|\mathbf{r}_q - \mathbf{u}\|_1 = \|\mathbf{r}_v\|_1 = \mathbf{1}^T \mathbf{r}_v$
  
- The PageRank equation is  $\mathbf{r}_v = c \mathbf{M} \mathbf{r}_v + (1-c) \mathbf{v}$
- Hence  $\mathbf{1}^T \mathbf{r}_v = c \underbrace{\mathbf{1}^T \mathbf{M}}_{\mathbf{1}^T} \mathbf{r}_v + (1-c) \mathbf{1}^T \mathbf{v}$  so that  $\mathbf{1}^T \mathbf{r}_v = \mathbf{1}^T \mathbf{v}$

As a result  $\|\mathbf{r}_q - \mathbf{u}\|_1 = \mathbf{1}^T \mathbf{v} < \sum \varepsilon d_j / D = \varepsilon$



# Scalability properties of Local PageRank using Approximate PageRank

(Francesco Barbato & Tommaso Boccatto, 2020)



Quasi-linear behaviour = **scalability** of Local PageRank



# Beware of the Lazy PageRank

which is suggested in the paper

□ Lazy PageRank  $r = a \underline{M_2} r + (1-a) q$

$$M_2 = b I + (1-b) M$$

□ **Lazy** because a fraction  $b$  of the times the surfer stays where she/he is

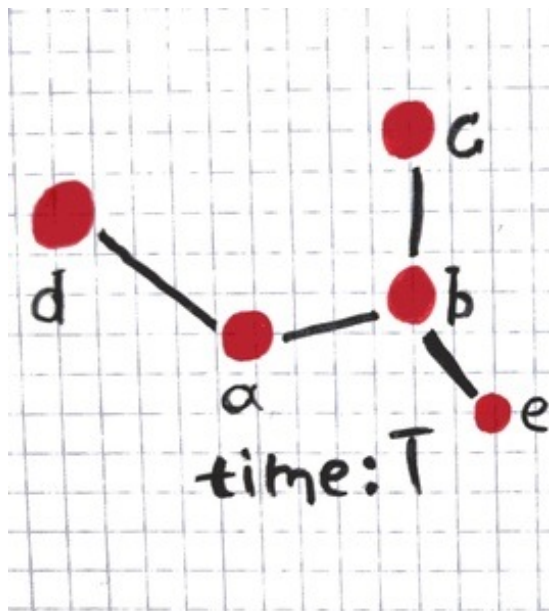
□ Equivalent to  $r = c \underline{M} r + (1-c) q$

$$c = a(1-b) / \underbrace{(1-ab)} < a$$

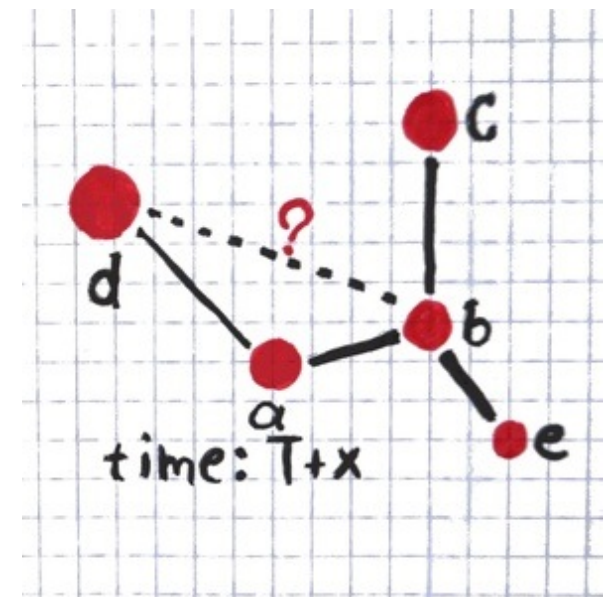
slower algorithm, as its convergence speed depends on  $a > c$ , better use  $c$  directly!

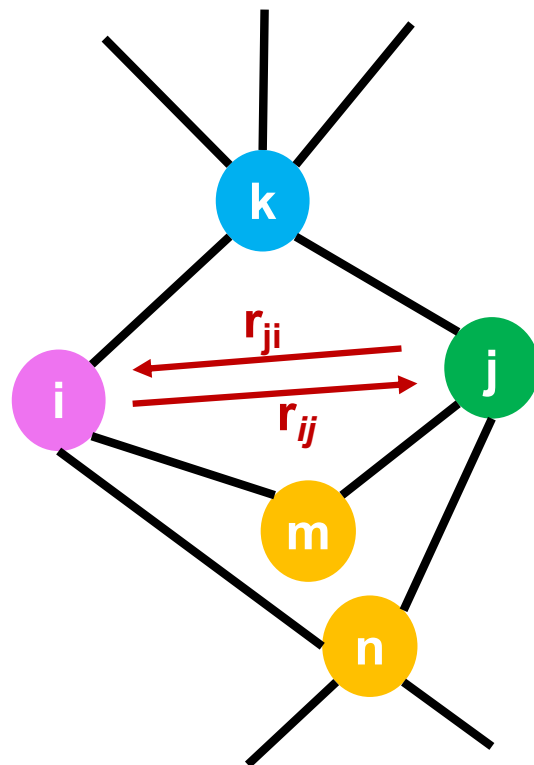


## Recommendation in social networks



Given a graph at time  $T$ , can we output a **ranked list** of links that are **predicted** to appear in the graph at time  $T+x$  ?





Local PageRank  
vector

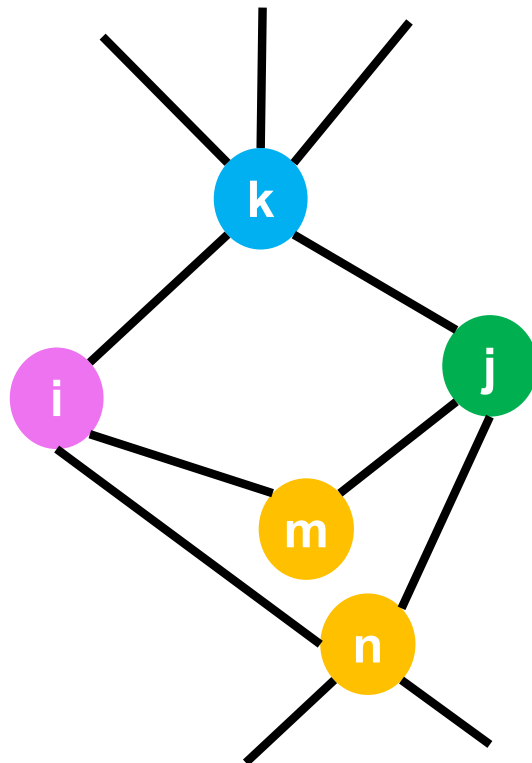
teleportation  
to node  $i$

$$\mathbf{r}_i = c \mathbf{M} \mathbf{r}_i + (1-c) \boldsymbol{\delta}_i$$

Likelihood of activating the link (i,j)

$$L_{\text{RWR}}(i,j) = r_{ij} + r_{ji}$$

Select the **highest** values of  $L_{\text{RWR}}$   
for **recommendation** purposes



$$L_{RA}(i,j) = \sum_{k \in N_i \cap N_j} 1/d_k$$

common  
neighbours

related to a two-hop RWR

$$\mathbf{r}_i \simeq (1-c) \sum_{n=0}^2 (c \mathbf{M})^n \boldsymbol{\delta}_i$$

to have

$$r_{ij} \simeq (1-c) c^2 / d_i L_{RA}(i,j)$$

$$L_{RWR}(i,j) \simeq (1-c) c^2 (1/d_i + 1/d_j) L_{RA}(i,j)$$



fraction of links correctly guessed  
(out of 100 recommendations)

<b>Precision</b>	CN	RA	LP	ACT
USAir	0.59	0.64	0.61	0.49
NetScience	0.26	0.54	0.30	0.19
Power	0.11	0.08	<b>0.13</b>	0.08
Yeast	0.67	0.49	0.68	0.57
C.elegans	0.12	0.13	<b>0.14</b>	0.07

RWR	HSM	LRW	SRW
0.65	0.28	0.64(3)	<b>0.67(3)</b>
<b>0.55</b>	0.25	0.54(2)	0.54(2)
0.09	0.00	0.08(2)	0.11(3)
0.52	0.84	<b>0.86(3)</b>	0.73(9)
0.13	0.08	<b>0.14(3)</b>	<b>0.14(3)</b>

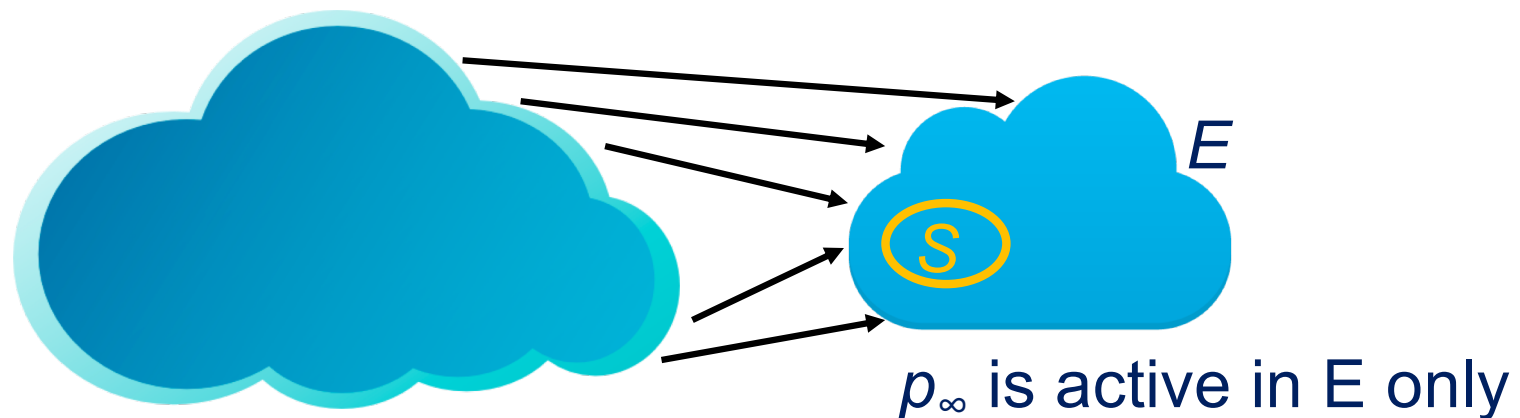
Among the best performance in **social networks**

But not strikingly good compared to simpler methods (e.g., RA = resource allocation)

- ❑ Bias the random walk towards a **topic specific teleport set  $S$**  of nodes, i.e., make sure that  $q$  is active in  $S$  only
- ❑  $S$  should contain only pages that are relevant to the topic

## Result

- ❑ The random walk **deterministically** ends in a small set  $E$ , containing  $S$ , and being in some sense close to it





**Tweet 1** is assigned  
to **Topic 1** !!!

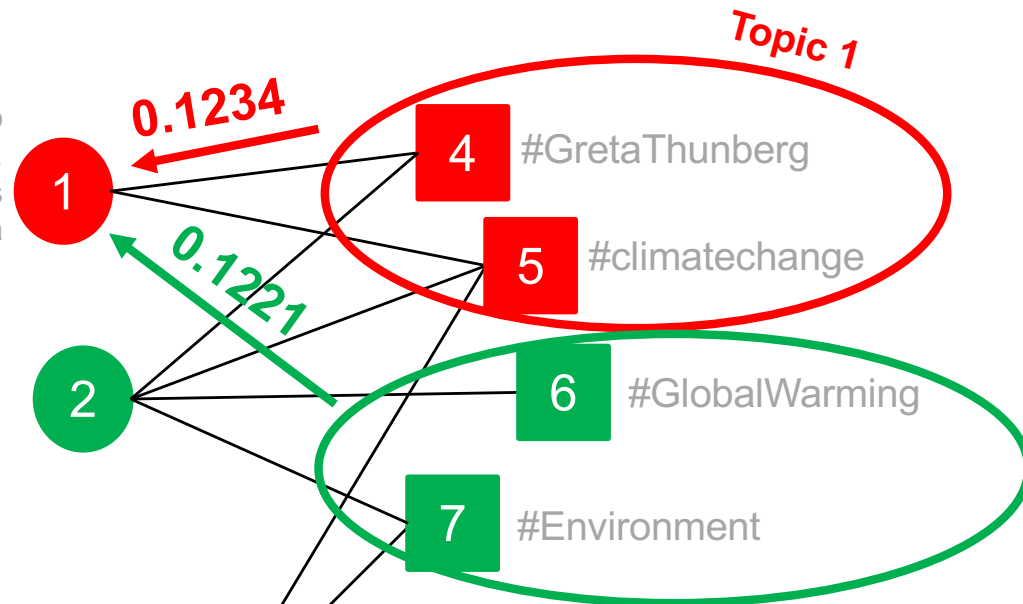
those who think they are crazy enough to  
change the world eventually do.  
#climatechange #ClimateCrisis  
#ClimateAction #GretaThunberg #Greta

Hopefully these kids will succeed where  
past generations have failed.  
#TheResistance #FBR #ClimateChange  
#Environment #GlobalWarming  
#GretaThunberg

The #environment can have a major effect  
on the human cardiovascular system. A  
new study has found an increase in heat-  
induced #heartattack risk in recent years.  
Could #ClimateChange be a risk factor?  
#longevity



Tweets



Hashtags

# Signed PageRank

modifications for signed networks



# PageRank in signed networks

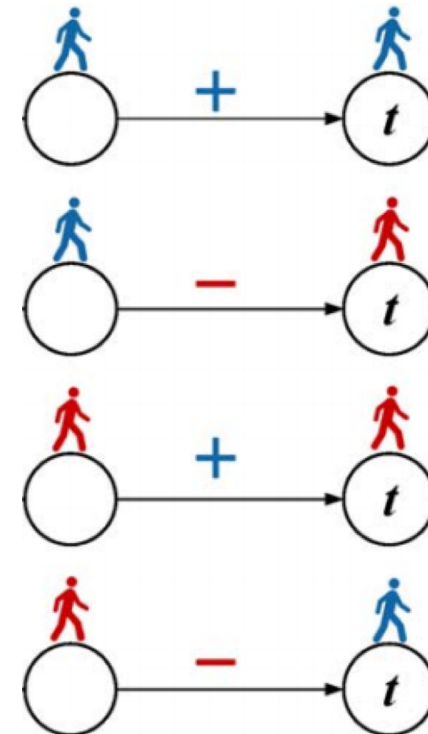
Jung, Jim, Sael, Kang, "Personalized ranking in signed networks using signed random walk with restart," 2016

<https://ieeexplore.ieee.org/iel7/7837023/7837813/07837935.pdf>

- Identify + (favourable) and – (adversarial) paths, i.e., ranking values  $r_+$  and  $r_-$  for positive and negative surfers
- Extract positive  $A_+$  and negative  $A_-$  contributions to  $A = A_+ - A_-$
- Normalize the absolute value, to get  $M_+$  and  $M_-$  (with normalized  $M_+ + M_-$ )
- Run a signed random walk

$$r_+ = c M_+ r_+ + c M_- r_- + (1-c) q$$

$$r_- = c M_- r_+ + c M_+ r_-$$







damping factor

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_+ & \mathbf{M}_- \\ \mathbf{M}_- & \mathbf{M}_+ \end{bmatrix} \text{ (column) normalized adjacency matrix}$$
$$\mathbf{r} = c \mathbf{M} \mathbf{r} + (1-c) \mathbf{q}_0$$

teleportation vector

$$\mathbf{q}_0 = \begin{bmatrix} \mathbf{q} \\ \mathbf{0} \end{bmatrix}$$

PageRank vector (centrality)  $\mathbf{r} = \begin{bmatrix} r_+ \\ r_- \end{bmatrix}$

signed centrality outcome  $\mathbf{r}_{+-} = r_+ - r_-$

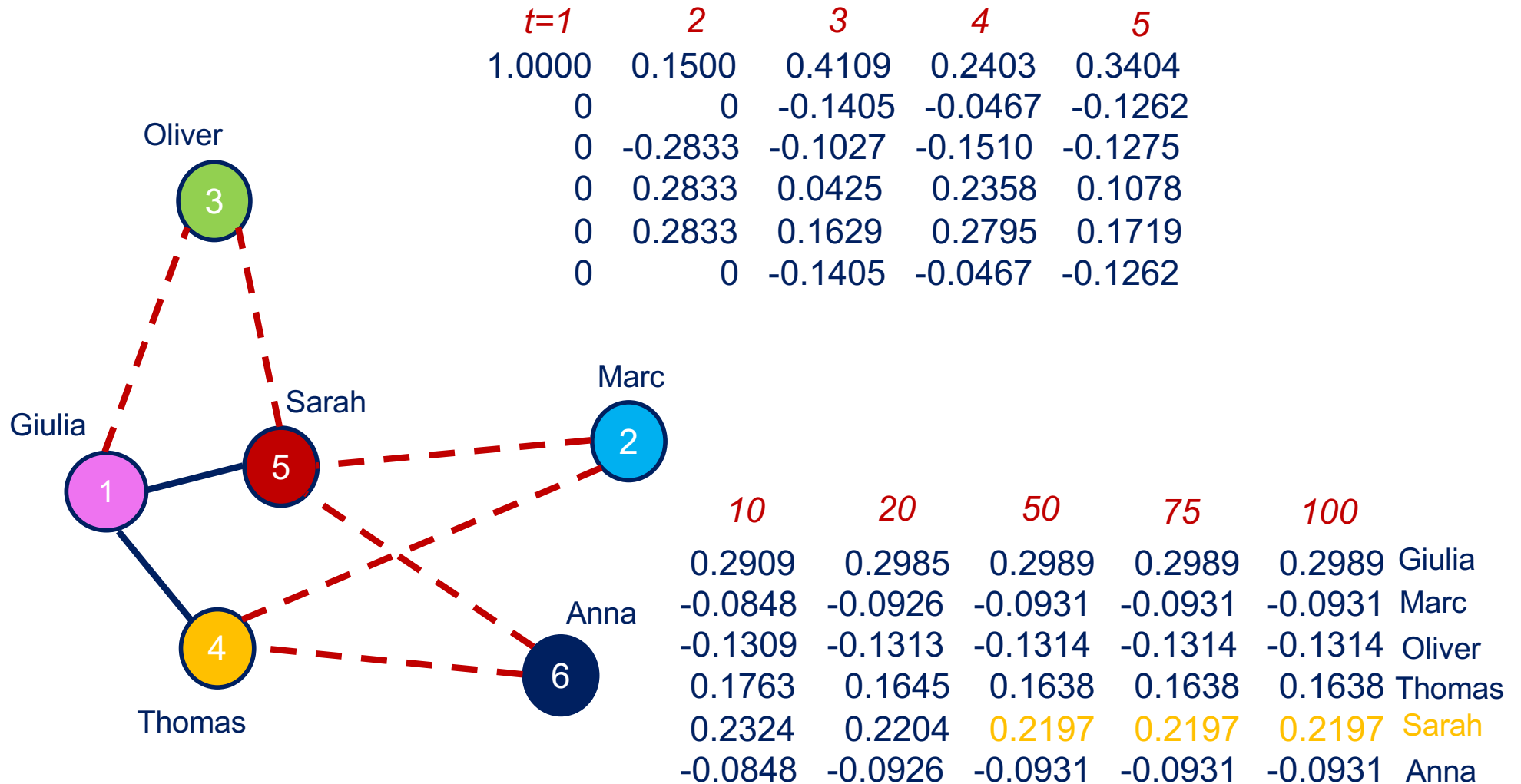
$$\mathbf{r}_{+-} = c \mathbf{M}_{+-} \mathbf{r}_{+-} + (1-c) \mathbf{q} \leftarrow \text{can be signed}$$

$$\mathbf{M}_{+-} = \mathbf{A} \text{diag}^{-1}(|\mathbf{A}|^T \mathbf{1})$$



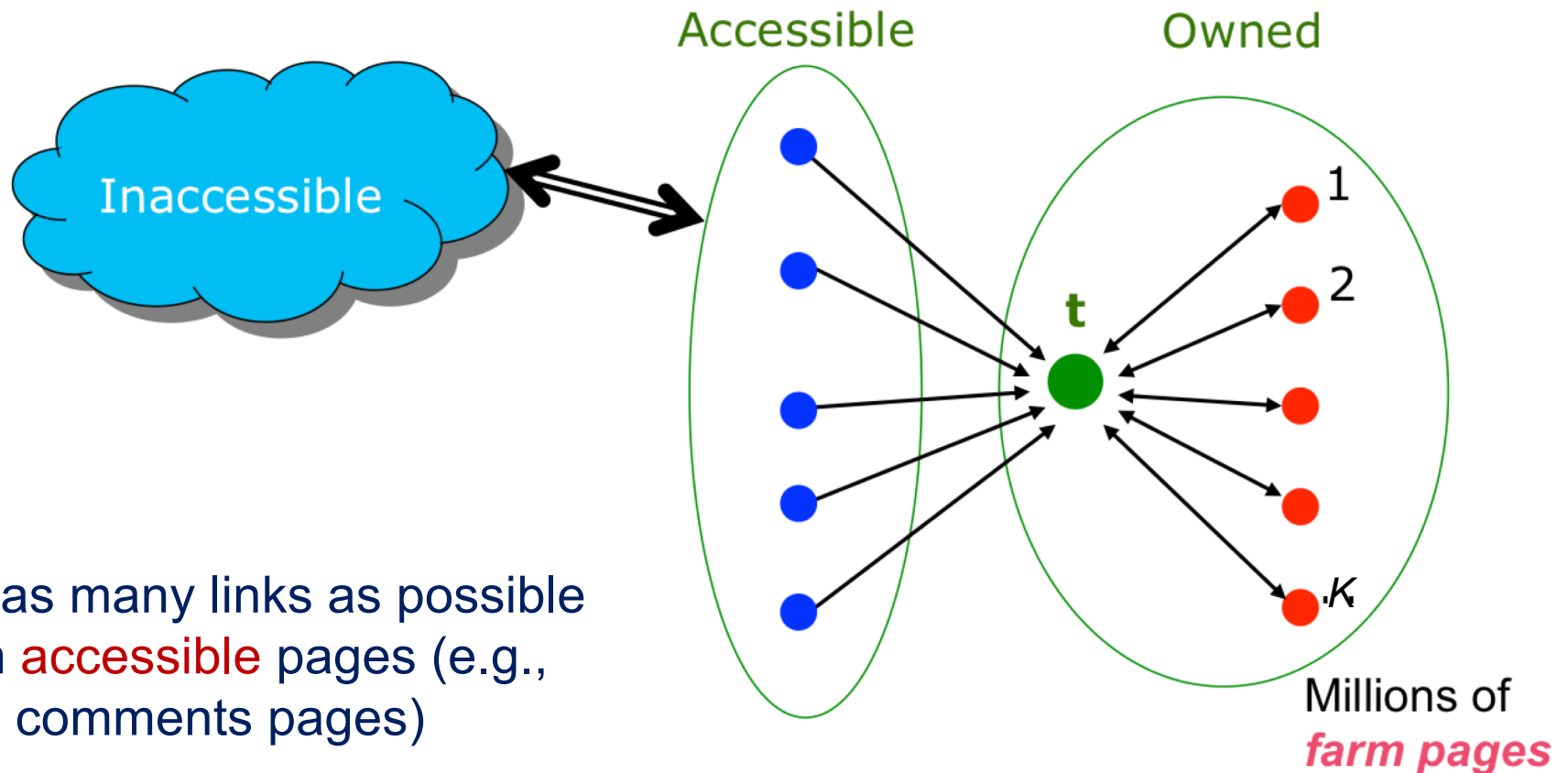
# Example

who's Giulia's best friend?



# Preventing spamming

on the role of the teleport vector



1. Get as many links as possible from **accessible** pages (e.g., blog comments pages)
2. Construct **link farm** to get a PageRank multiplier effect



## Web

Results 1 - 10 of about 969,000 for [miserable failure](#). (0.06 seconds)

### [Biography of President George W. Bush](#)

Biography of the president from the official White House web site.

[www.whitehouse.gov/president/gwbbio.html](http://www.whitehouse.gov/president/gwbbio.html) - 29k - [Cached](#) - [Similar pages](#)

[Past Presidents](#) - [Kids Only](#) - [Current News](#) - [President](#)

[More results from www.whitehouse.gov »](#)

### [Welcome to MichaelMoore.com!](#)

Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ...

[www.michaelmoore.com/](http://www.michaelmoore.com/) - 35k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)

### [BBC NEWS | Americas | 'Miserable failure' links to Bush](#)

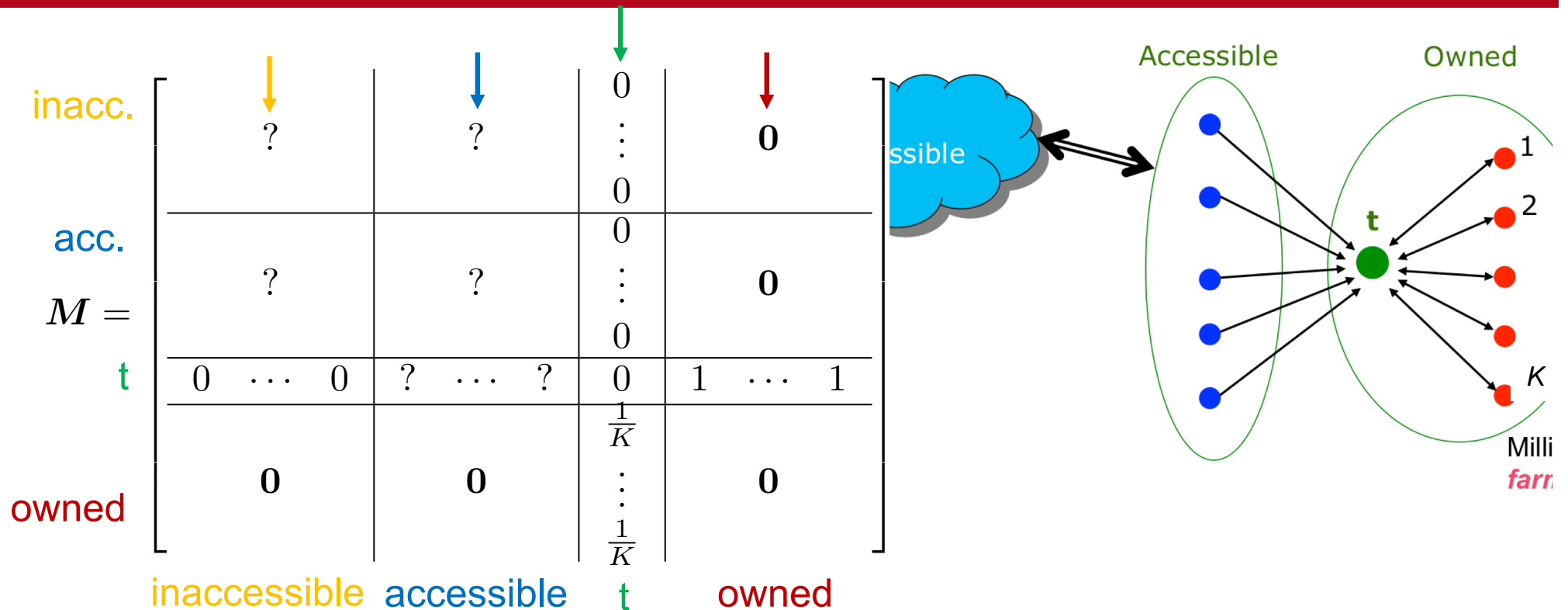
Web users manipulate a popular search engine so an unflattering description leads to the president's page.

[news.bbc.co.uk/2/hi/americas/3298443.stm](http://news.bbc.co.uk/2/hi/americas/3298443.stm) - 31k - [Cached](#) - [Similar pages](#)

### [Google's \(and Inktomi's\) Miserable Failure](#)

A search for **miserable failure** on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ...

[searchenginewatch.com/sereport/article.php/3296101](http://searchenginewatch.com/sereport/article.php/3296101) - 45k - Sep 1, 2005 - [Cached](#) - [Similar pages](#)



$$r = c M r + (1 - c) q$$



$$r_t = a + cK r_o + (1 - c) q_o$$

$$r_o = c \frac{1}{K} r_t + (1 - c) q_o$$

ranking due to accessible pages
teleportation value to pages owned by the spammer



ranking due to  
**accessible** pages

teleportation value to pages  
owned by the spammer

$$r_t = \frac{a}{1 - c^2} + \frac{cK + 1}{1 + c} q_o$$

scaling factor ( $\approx 3.6$ )

**spam** factor (can be made  
as large as desired)

## solution

teleport only to **trusted** pages (i.e., set  $q_o = 0$ )  
can also be used as a method to **identify** spam farms

# Row-normalized PageRank

For spreading information over the network





PageRank equation

$$\mathbf{r} = c \mathbf{M} \mathbf{r} + (1-c) \mathbf{q}$$

row-normalized

$$\mathbf{M} = \text{diag}^{-1}(\mathbf{d}) \mathbf{A}, \mathbf{d} = \mathbf{A} \mathbf{1} \quad \mathbf{M} \mathbf{1} = \mathbf{1}$$

Markov chain

$$\mathbf{p}_{t+1} = c \mathbf{M} \mathbf{p}_t + (1-c) \mathbf{q}$$

$$\mathbf{p}_0 = \mathbf{q}$$

$$\mathbf{M}_1 = c \mathbf{M} + (1-c) \mathbf{q} \mathbf{v}^T$$

$$\mathbf{v}^T \mathbf{M} = \mathbf{v}^T$$

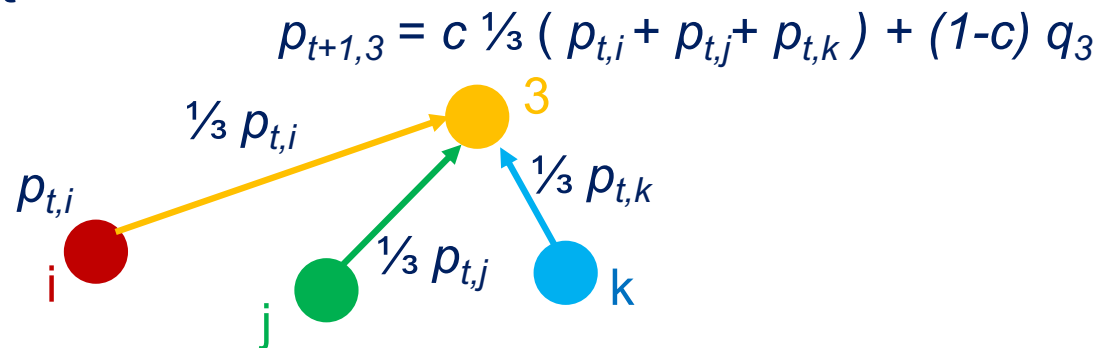
$$\mathbf{v}^T \mathbf{q} = 1$$

same properties of column-normalized PageRank:

- ❑  $\mathbf{M}_1$  has **one** eigenvalue equal to 1
- ❑ The remaining eigenvalues satisfy  $|\lambda| \leq c$



A node gathers the average value of the neighbour nodes pointing to it



It is a way of **spreading** the original information  $q$  over the network





- ❑ This is the metric to be used it for resizing nodes according to their importance
- ❑ Provides elaborate information on the **relevance** of nodes in the network
- ❑ For directed networks, it can be used in both its **authority** and **hub** forms
- ❑ Can also be put in the form of a PageRank **distribution**
- ❑ Can be used in different useful ways, e.g., to evaluate **similarity** or closeness, to **spread** information
- ❑ Exploit its potential at your best

# HITS centrality

a (less interesting) alternative to PageRank



## HITS – hubs and authorities

Kleinberg, J.M.

1999

«Authoritative sources in a  
hyperlinked environment»

*Journal of the ACM*

<https://www.cs.cornell.edu/home/kleinber/auth.pdf>

Conceptually similar to PageRank

Provides scores for **authorities** and **hubs**,  
separately, as PageRank can do

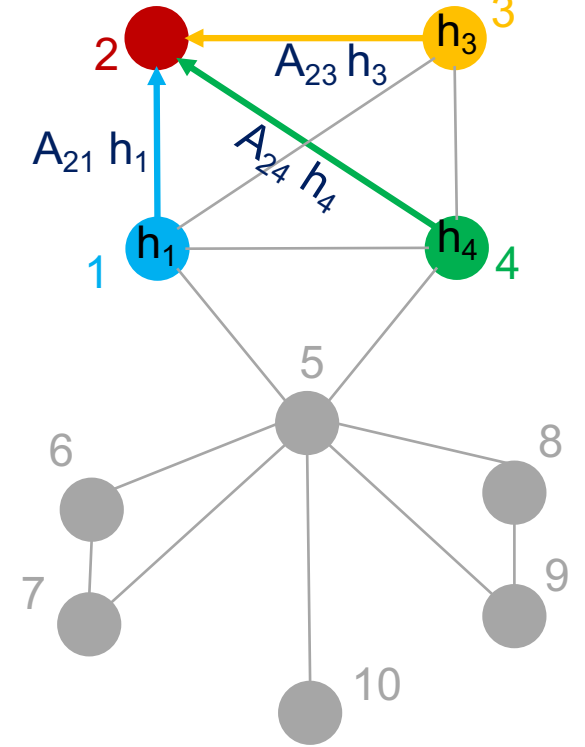
We **deprecate** its use



$A_{2,4}$  = weight of connection 4  $\rightarrow$  2

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ a_{10} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \\ h_6 \\ h_7 \\ h_8 \\ h_9 \\ h_{10} \end{bmatrix}$$

$$\begin{aligned} a_2 &= A_{21} h_1 + A_{23} h_3 + A_{24} h_4 \\ &= h_1 + h_3 + h_4 \end{aligned}$$



$$\underset{\substack{\uparrow \\ \text{authority scores}}}{a} = A \underset{\substack{\uparrow \\ \text{hub scores}}}{h}$$

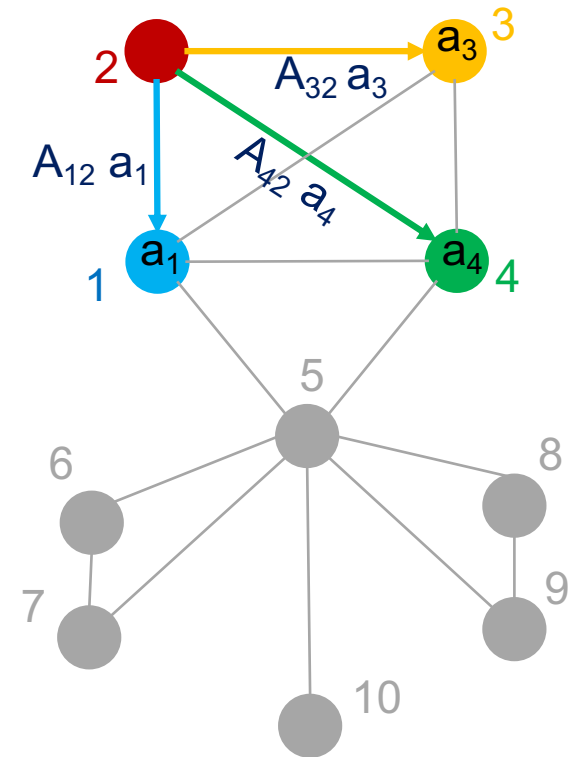


$A_{3,2}$  = weight of connection 2  $\rightarrow$  3

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \\ h_6 \\ h_7 \\ h_8 \\ h_9 \\ h_{10} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}^T \cdot \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ a_{10} \end{bmatrix}$$

$$h = A^T a$$

$$h_2 = A_{12} a_1 + A_{32} a_3 + A_{42} a_4$$
$$= a_1 + a_3 + a_4$$







$$a = c_a \cdot Ah$$
$$h = c_h \cdot A^T a$$

hubs

$$h = c M h$$

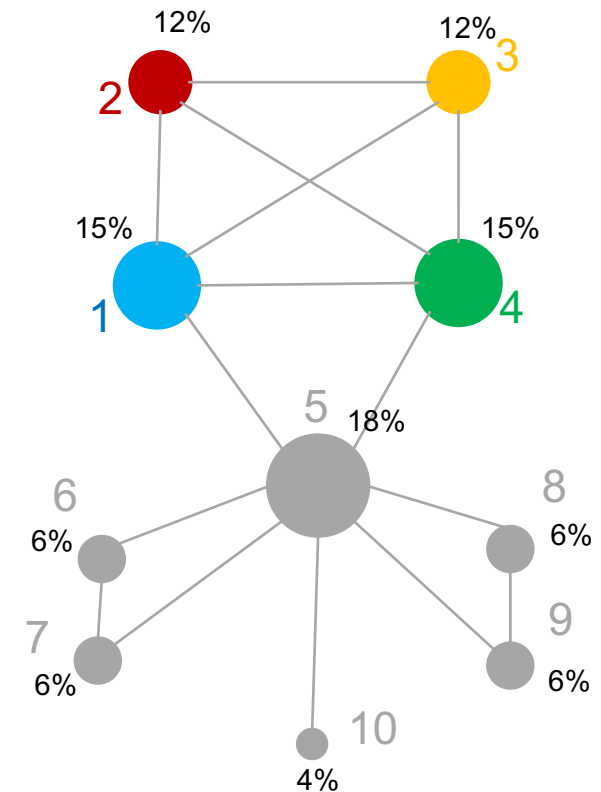
$$M = A^T A$$

$$c = c_a c_h$$

authorities

$$a = c_a \cdot Ah$$

- ❑ The formula says we are interested in the (principal) **eigenvector** of matrix  $M = A^T \cdot A$
- ❑ Can be obtained by standard linear algebra algorithms





0. Start from an  
initial guess  $\mathbf{a}_0$

1. Let the time go by

$$\mathbf{a}_{t+1} = \mathbf{M} \mathbf{a}_t$$

product by a sparse  
matrix (twice)  $\mathbf{M} = \mathbf{A} \mathbf{A}^T$

2. Keep normalizing  
(divide  $\mathbf{a}_{t+1}$  by the sum  
of elements)

3. Stop when  $\mathbf{a}$   
converges (few iterations?)



- ❑  $\|\mathbf{a}_t - \mathbf{a}_\infty\|_2 \leq \sqrt{N} \cdot (\lambda_2/\lambda_1)^t$
- ❑  $\lambda_1$  largest **eigenvalue** of  $\mathbf{M}$
- ❑  $\lambda_2$  second largest eigenvalue of  $\mathbf{M}$
- ❑ Triang. inequality  $\|\mathbf{a}_t - \mathbf{a}_{t+1}\|_2 \leq 2\sqrt{N} \cdot (\lambda_2/\lambda_1)^t$

Worst case result:

- ❑ **Precision**  $\varepsilon$  implies:  $\|\mathbf{a}_t - \mathbf{a}_{t+1}\|_2 < \varepsilon$

- ❑ **Iterations** required:  $t = \lceil [\ln(2/\varepsilon) + \frac{1}{2}\ln(N)] / \ln(\lambda_1/\lambda_2) \rceil$

$N = 10^9 \rightarrow 10.3$

slow if  $\lambda_2$  close to  $\lambda_1$


$10^{-3}$  precision  $\rightarrow 7.6$

# Eigenvector and Katz centralities

other (less interesting) alternatives to PageRank



	with constant term	without constant term
normalized	PageRank $\mathbf{r} = c \mathbf{M} \mathbf{r} + (1-c) \mathbf{q}$	Degree $\mathbf{r} = \mathbf{M} \mathbf{r}$
unnormalized	Katz $\mathbf{r} = c \mathbf{A} \mathbf{r} + \mathbf{1}$	Eigenvector $\mathbf{r} = c \mathbf{A} \mathbf{r}$


$$\begin{aligned} \mathbf{r} &= (\mathbf{I} - c \mathbf{A})^{-1} \mathbf{1} \\ &= \sum (c \mathbf{A})^k \mathbf{1} \end{aligned}$$

The absence of normalization makes them less robust and meaningful compared to PageRank

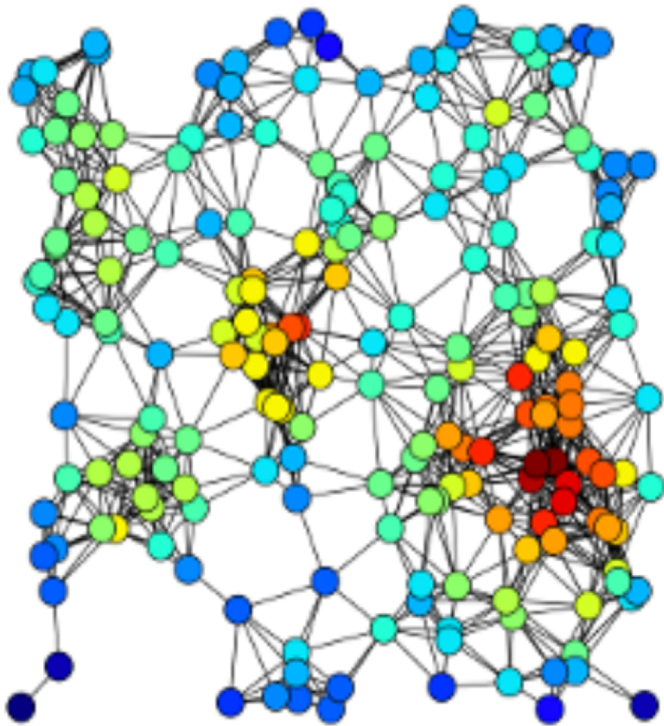
They are **deprecated**



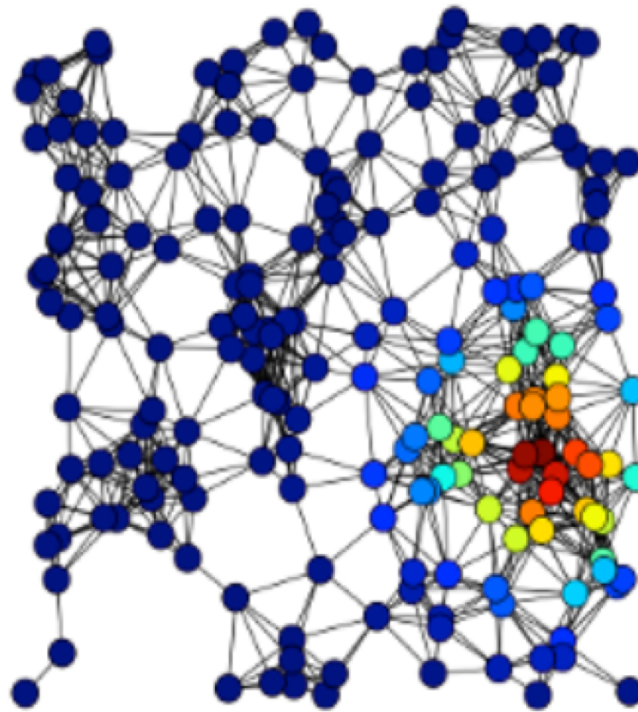
# Eigenvector and Katz centralities

their graphical interpretation

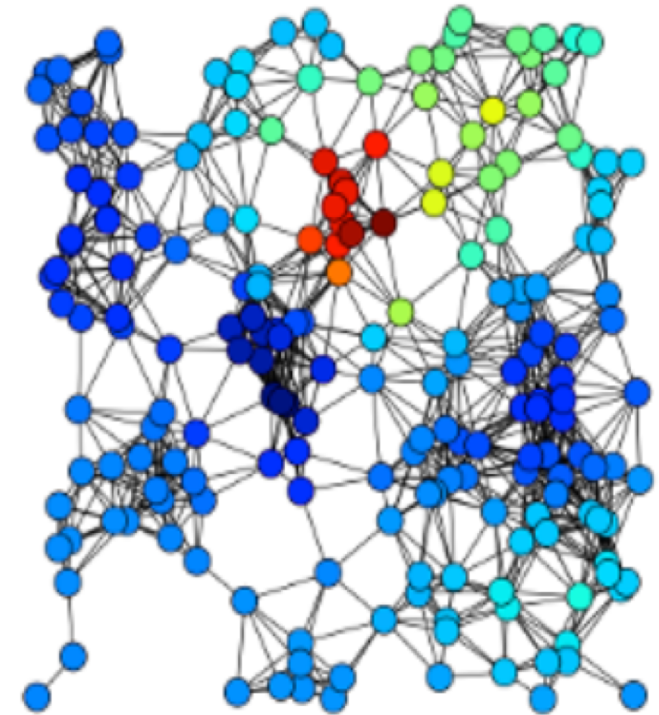
Degree



Eigenvector



Katz



# Closeness and Harmonic centralities

importance of nodes as spreaders of information



## Closeness centrality

---

From Wikipedia, the free encyclopedia



In a **connected graph**, **closeness centrality** (or **closeness**) of a node is a measure of **centrality** in a **network**, calculated as the reciprocal of the sum of the length of the **shortest paths** between the node and all other nodes in the graph. Thus, the more central a node is, the *closer* it is to all other nodes.

Closeness was defined by Bavelas (1950) as the **reciprocal** of the **farness**,<sup>[1][2]</sup> that is:

$$C(x) = \frac{1}{\sum_y d(y, x)}.$$

where  $d(y, x)$  is the **distance** between vertices  $x$  and  $y$ . When

**Rationale:** the node which is the easiest to reach, the one which is the best for spreading information

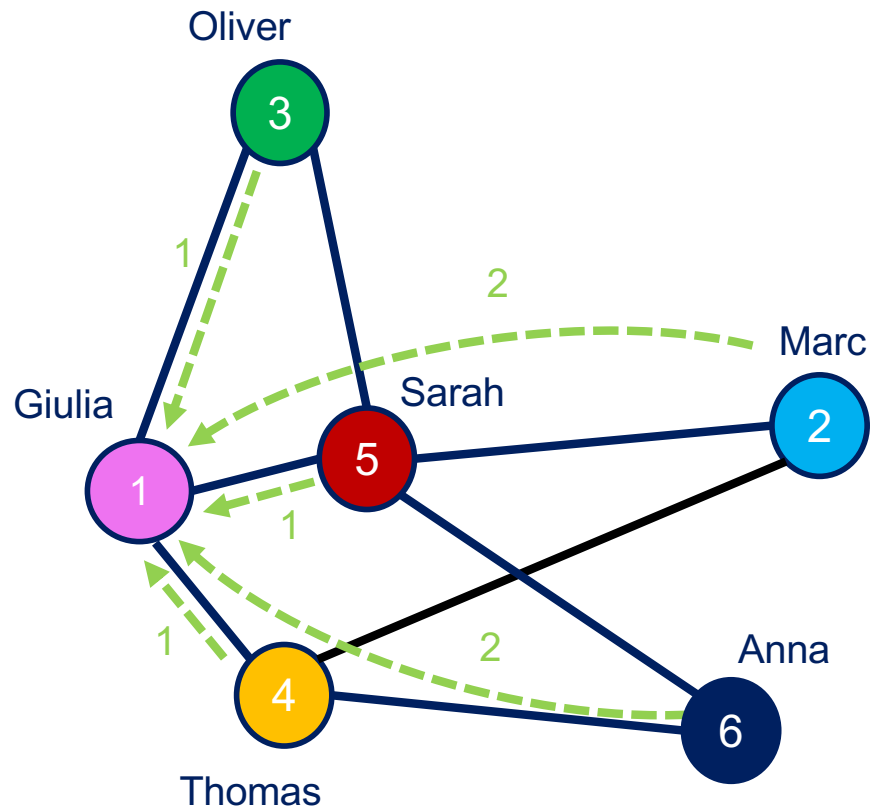




# An example

on how to calculate closeness centrality

count the lengths of the shortest paths  
leading to Giulia  
 $1 + 2 + 1 + 2 + 1 = 7$



## Closeness

0.1429 Giulia  
0.1250 Marc  
0.1250 Oliver  
0.1429 Thomas  
0.1667 Sarah  
0.1250 Anna

Sarah is the preferred node for spreading information

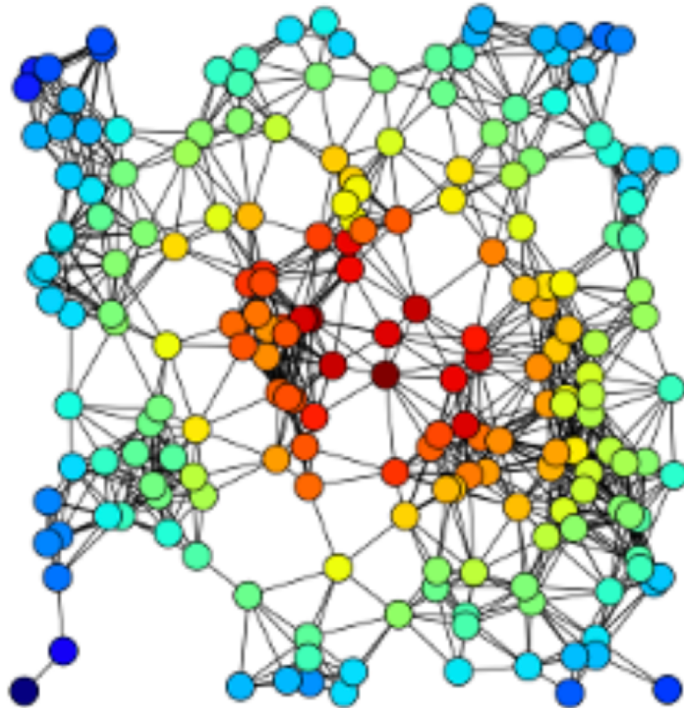
$$C(\text{Giulia}) = 1/7 \\ = 0.1429$$



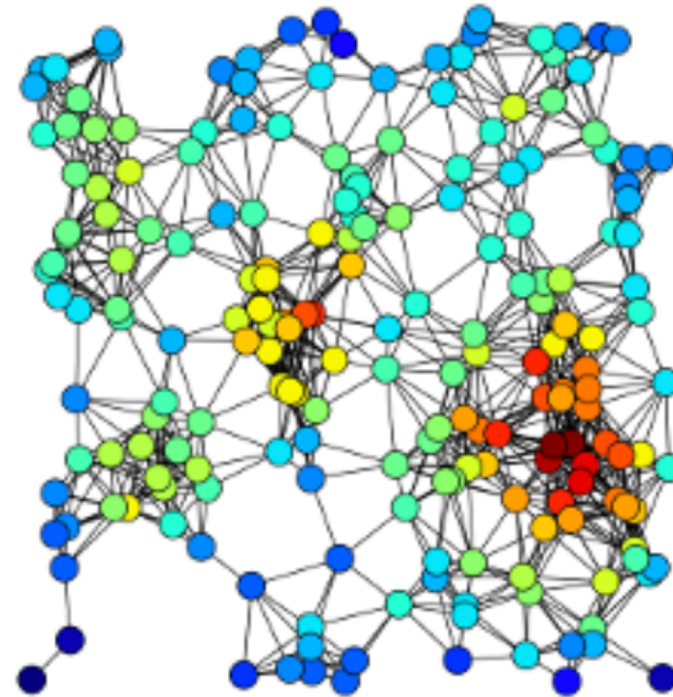
# Closeness versus degree centrality

a graphical interpretation

Closeness



Degree





## In disconnected graphs [\[ edit \]](#)

---

When a graph is not [strongly connected](#), a widespread idea is that of using the sum of reciprocal of distances, instead of the reciprocal of the sum of distances, with the convention  $1/\infty = 0$ :

$$H(x) = \sum_{y \neq x} \frac{1}{d(y, x)}.$$

The most natural modification of Bavelas's definition of closeness is following the general principle proposed by [Marchiori and Latora \(2000\)<sup>\[3\]</sup>](#) that in graphs with infinite distances the harmonic mean behaves better than the arithmetic mean. Indeed, Bavelas's closeness can be described as the denormalized reciprocal of the [arithmetic mean](#) of distances, whereas harmonic centrality is the denormalized reciprocal of the [harmonic mean](#) of distances.

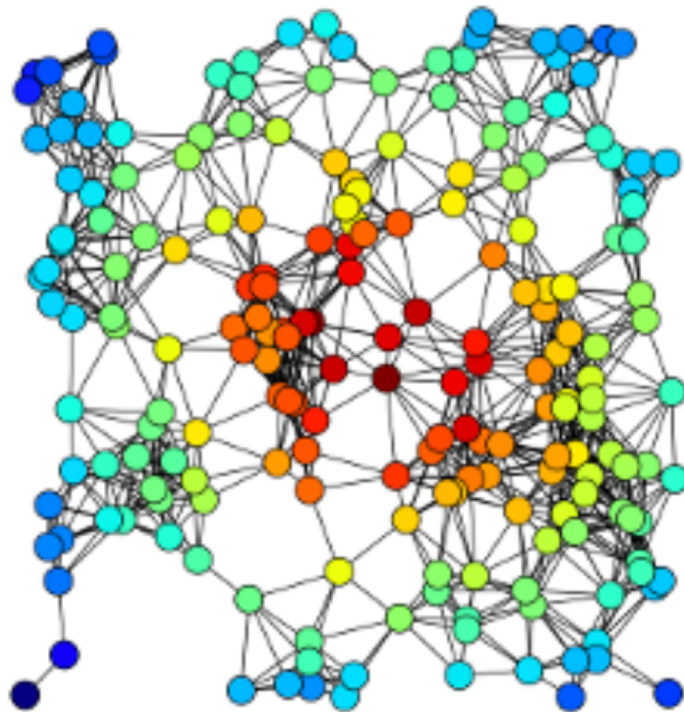




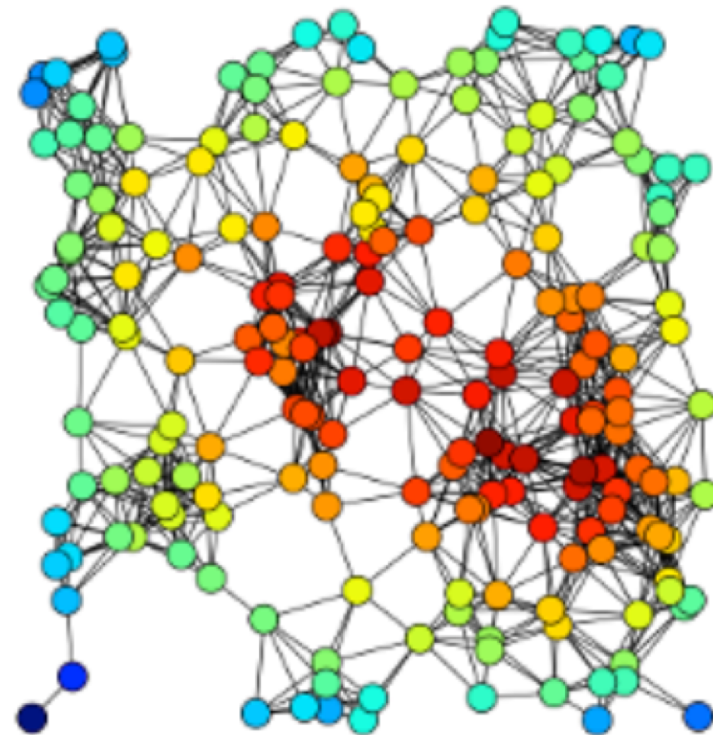
# Closeness versus harmonic centrality

a graphical interpretation

Closeness



Harmonic



# Betweenness centrality

importance of nodes as bridges or brokers



## Betweenness centrality

---

From Wikipedia, the free encyclopedia



In [graph theory](#), **betweenness centrality** is a measure of [centrality](#) in a [graph](#) based on [shortest paths](#). For every pair of vertices in a connected graph, there exists at least one shortest path between the vertices such that either the number of edges that the path passes through (for unweighted graphs) or the sum of the weights of the edges (for weighted graphs) is minimized. The betweenness centrality for each [vertex](#) is the number of these shortest paths that pass through the vertex.

Betweenness centrality was devised as a general measure of centrality:<sup>[1]</sup> it applies to a wide range of problems in network theory, including problems related to social [networks](#), biology, transport and scientific cooperation. Although earlier authors have intuitively described centrality as based on betweenness, [Freeman \(1977\)](#) gave the first formal definition of betweenness centrality.

*Rationale: the node which takes  
you elsewhere  
(bridge, broker)*



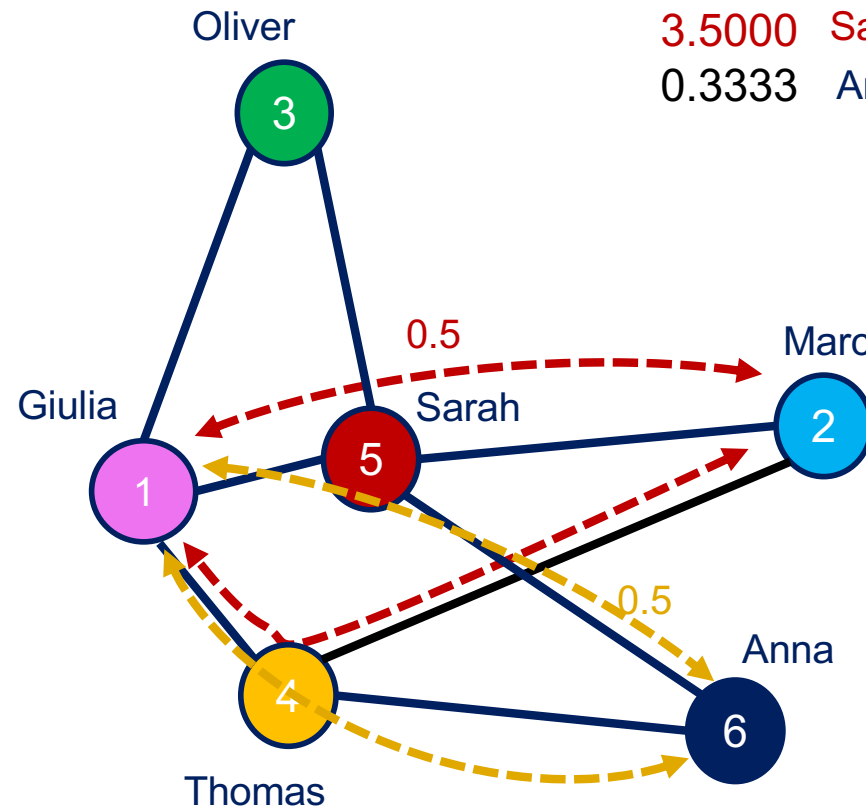
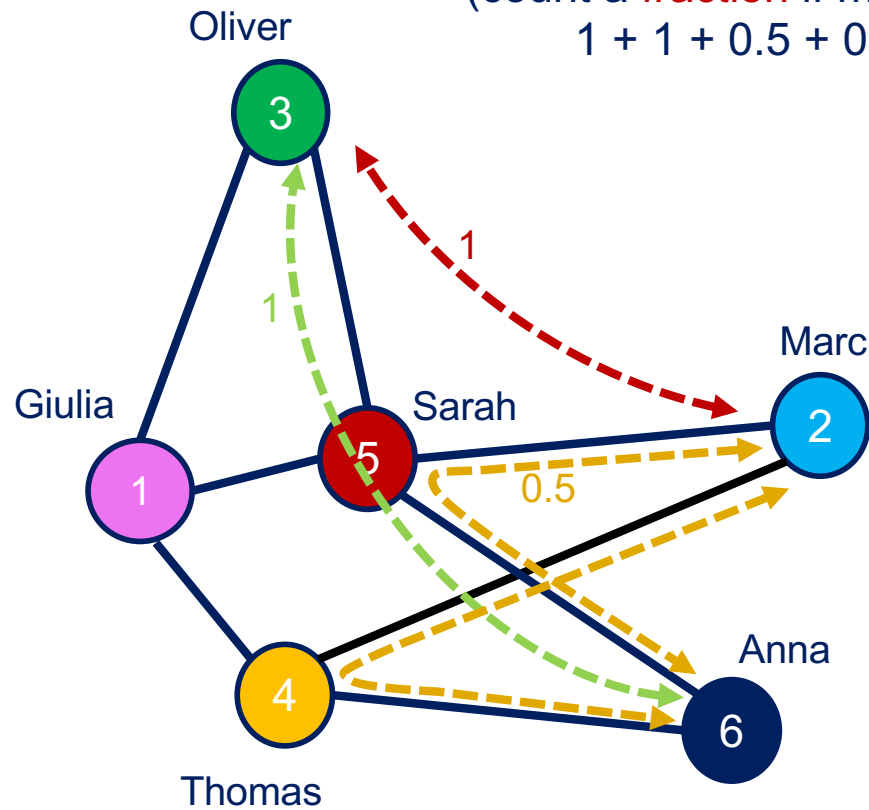
# An example

on how to calculate betweenness centrality

count the # of shortest paths  
passing through Sarah  
(count a **fraction** if more than one path)  
 $1 + 1 + 0.5 + 0.5 + 0.5 = 3.5$

## Betweenness

1.3333	Giulia
0.3333	Marc
0	Oliver
1.5000	Thomas
3.5000	Sarah
0.3333	Anna

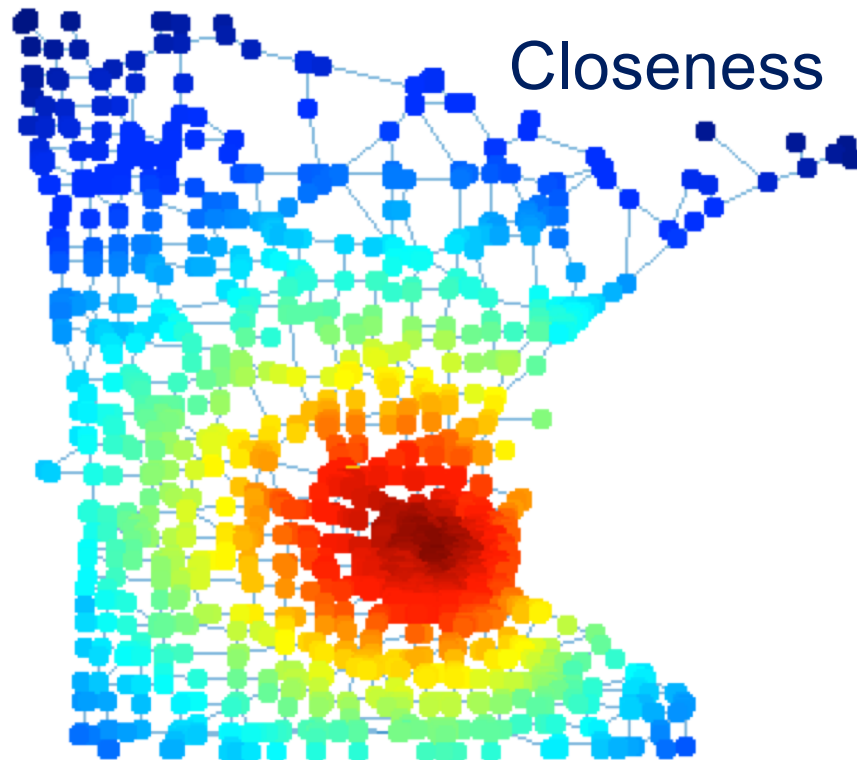




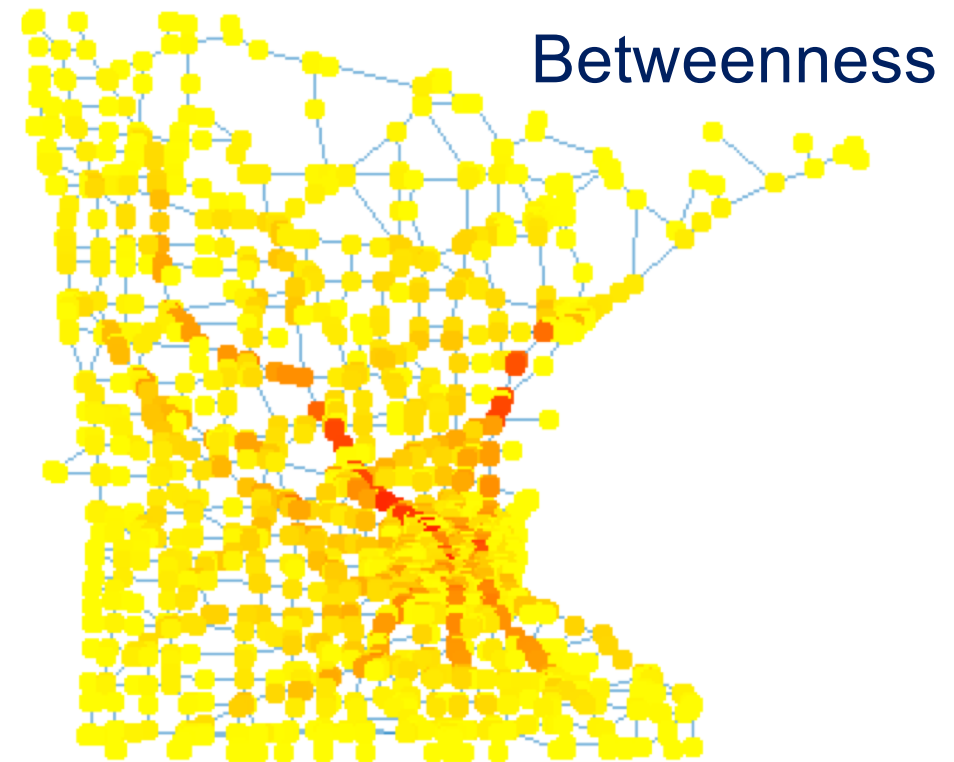
# Closeness vs betweenness centrality

a graphical interpretation

Minnesota road network



**Closeness** is a measure of **center of gravity** (best node to spread info)



**Betweenness** is a measure of **brokerage** (i.e., being a bridge)

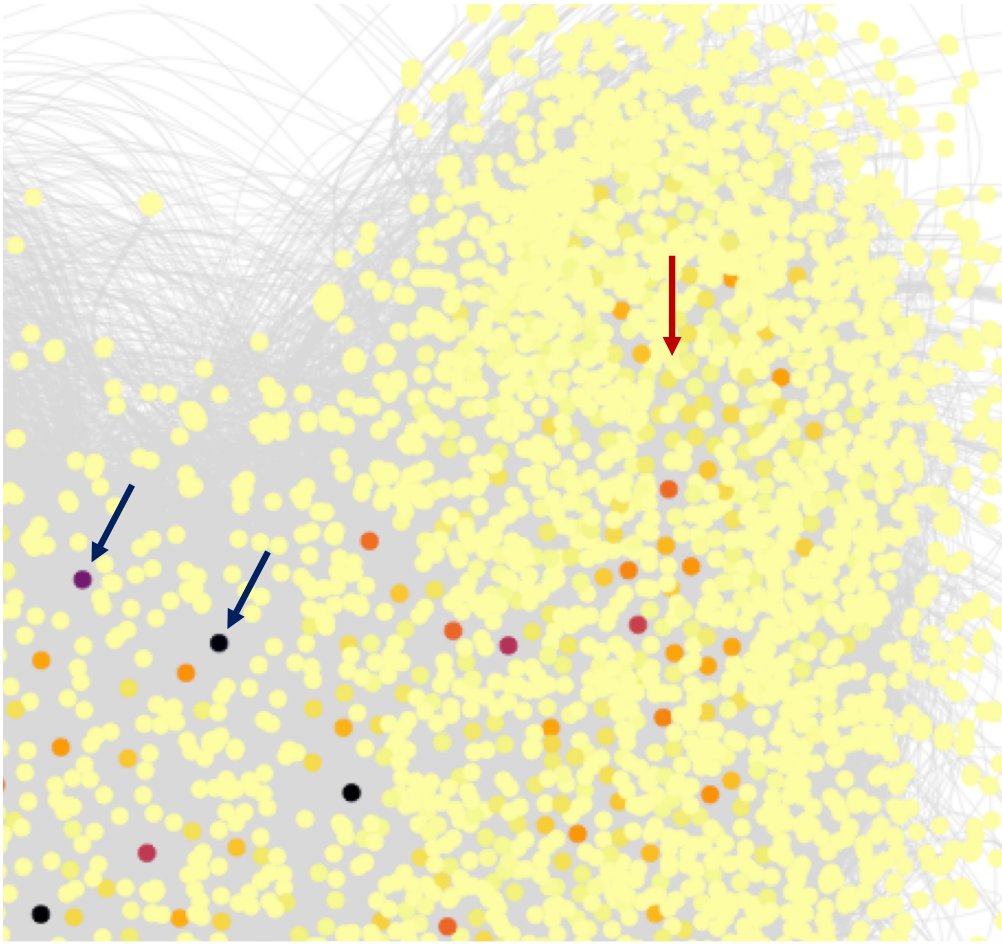




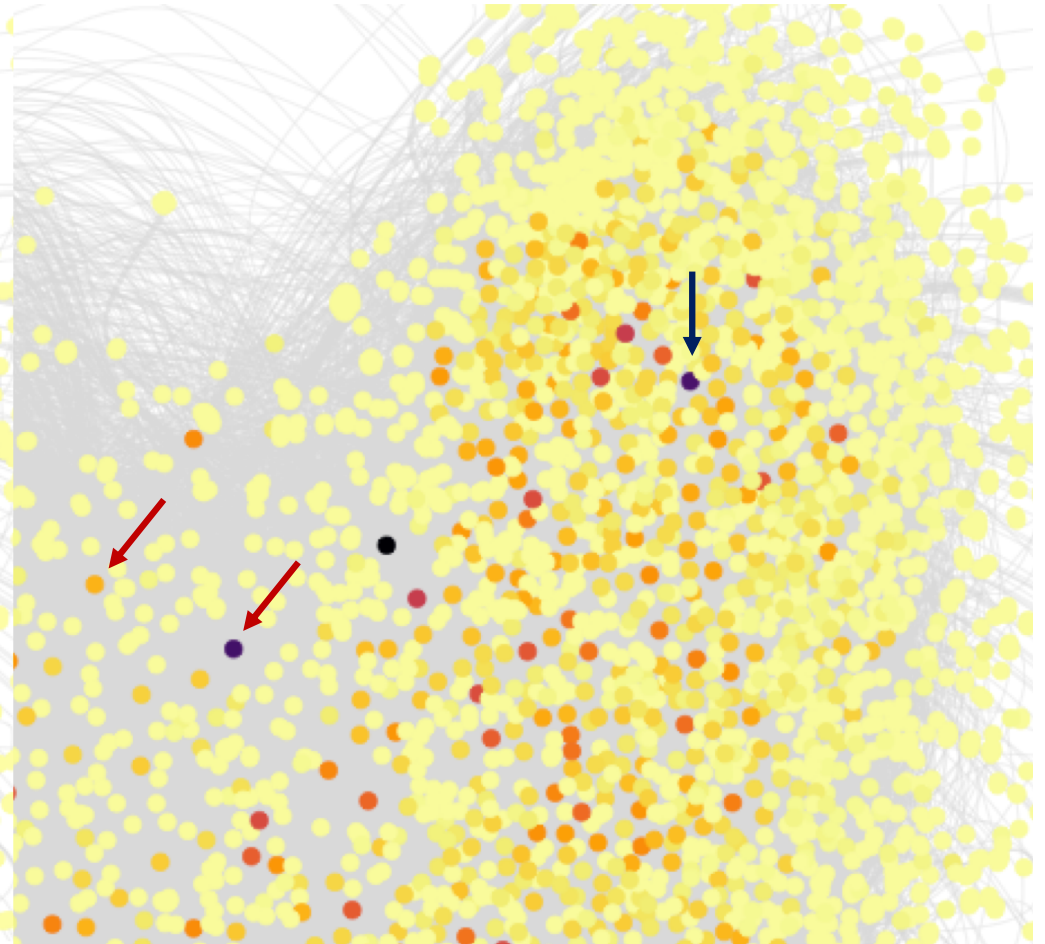
# Betweenness vs PageRank centrality

wiki vote network

## Betweenness



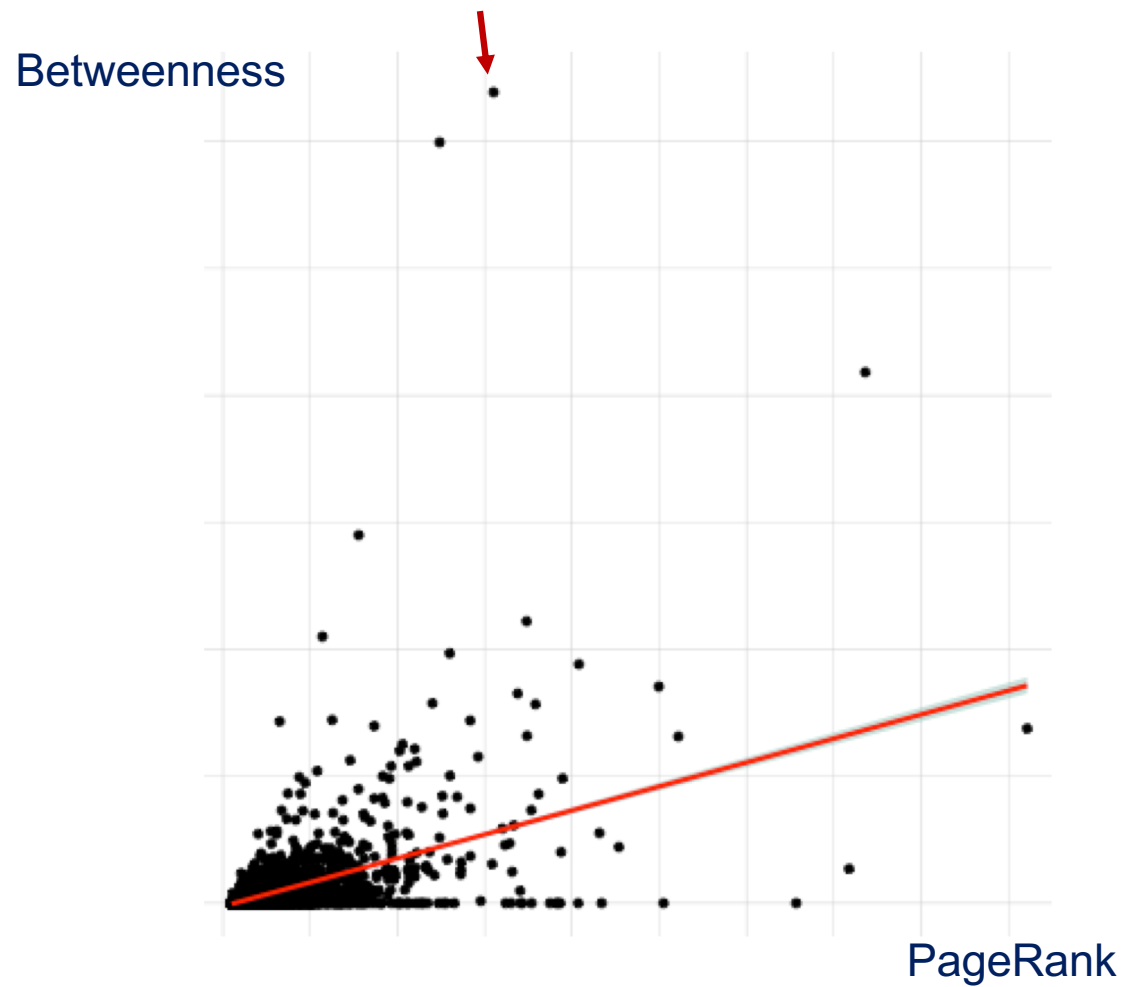
## PageRank





# Betweenness vs PageRank centrality

a correlation view



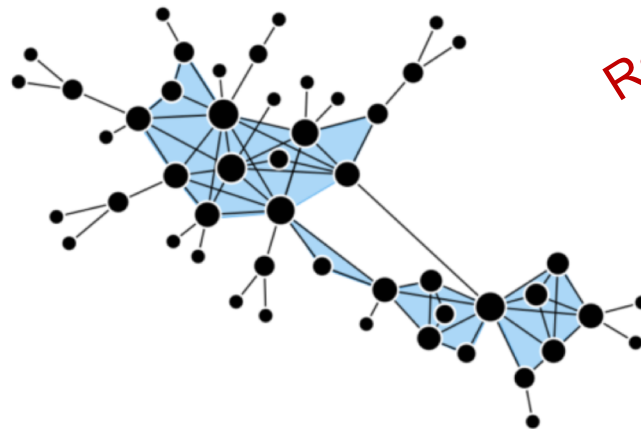
# Clustering coefficient

how tightly linked is the network locally

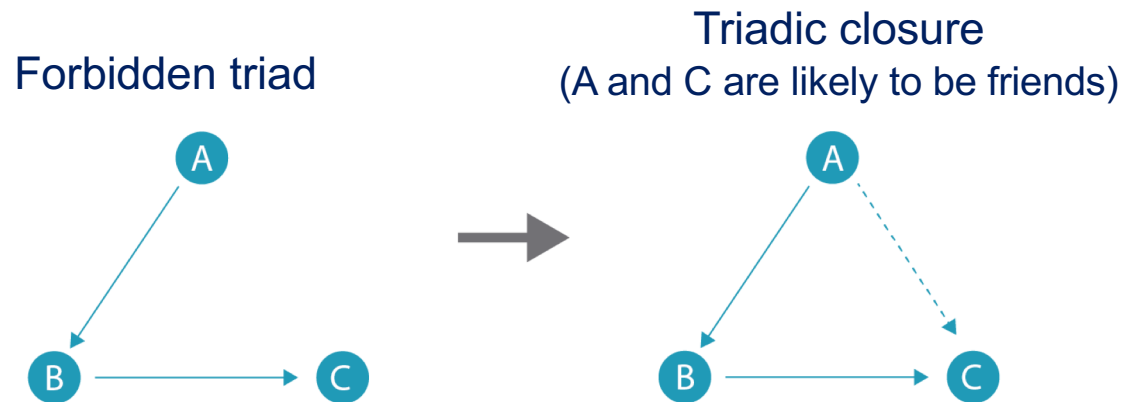


## Local clustering coefficient [\[ edit \]](#)

The **local clustering coefficient** of a **vertex** (node) in a **graph** quantifies how close its **neighbours** are to being a **clique** (complete graph). **Duncan J. Watts** and **Steven Strogatz** introduced the measure in 1998 to determine whether a graph is a **small-world network**.



**Rationale:** how strongly connected is the network locally / general indication of the graph's tendency to be organized into clusters



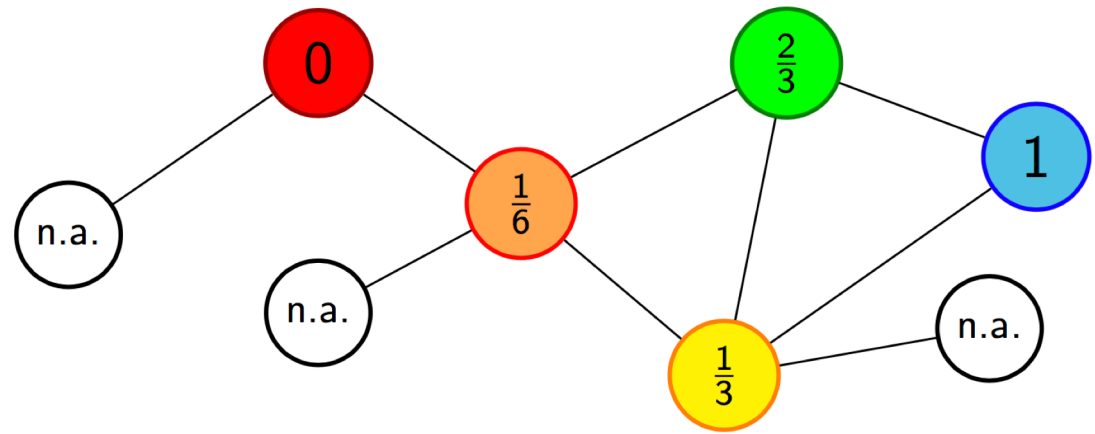
## Triadic closure

- ❑ A and C are likely to have the opportunity to meet because they have a common friend B
- ❑ The fact that A and C is friends with B gives them the basis of **trusting** each other
- ❑ B may have the **incentive** to bring A and C together, as it may be hard for B to maintain disjoint relationships



# Local clustering coefficient

a measure of triadic closures



**Local Clustering** coefficient  $C_i$  counts the **fraction** of pairs of neighbours  $N_i$  which form a triadic closure with node  $i$

$$C_i = \frac{1}{|\mathcal{N}_i|(|\mathcal{N}_i| - 1)} \sum_{\substack{(j,k) \in \mathcal{N}_i^2 \\ i \neq k}} tc_{i,j,k}$$

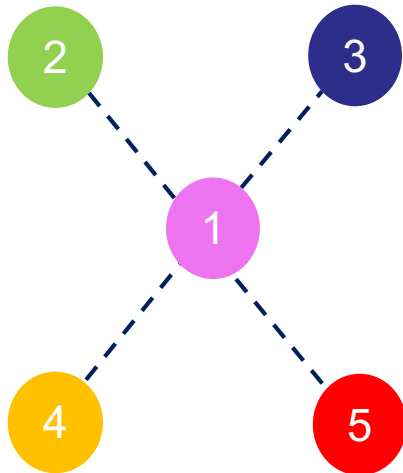
equal to  $\text{diag}(\mathbf{A}^3)$

where  $tc_{ijk} = 1$  if the triplet  $(i,j,k)$  forms a triadic closure, and zero otherwise



not connected  
neighbourhood

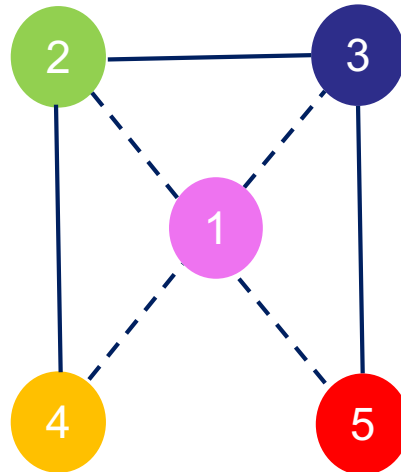
$$\langle C \rangle = 0$$



$$C_1 = 0$$

weakly connected  
neighbourhood

$$\langle C \rangle = 0.766$$



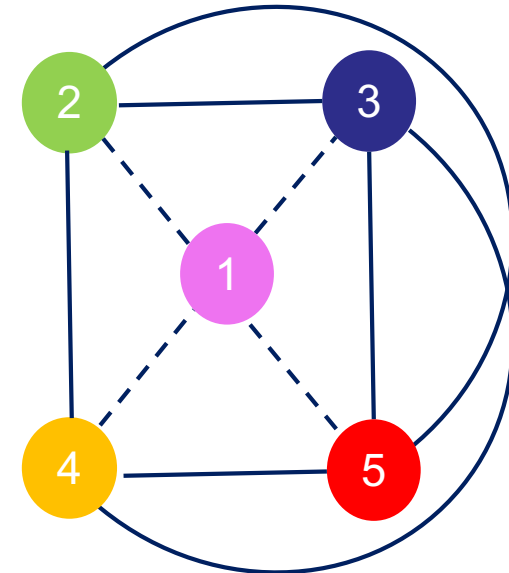
$$C_1 = \frac{1}{2} = \frac{3}{(4 \times 3/2)}$$

$$C_2 = C_3 = \frac{2}{3}$$

$$C_4 = C_5 = 1$$

strongly connected  
neighbourhood

$$\langle C \rangle = 1$$



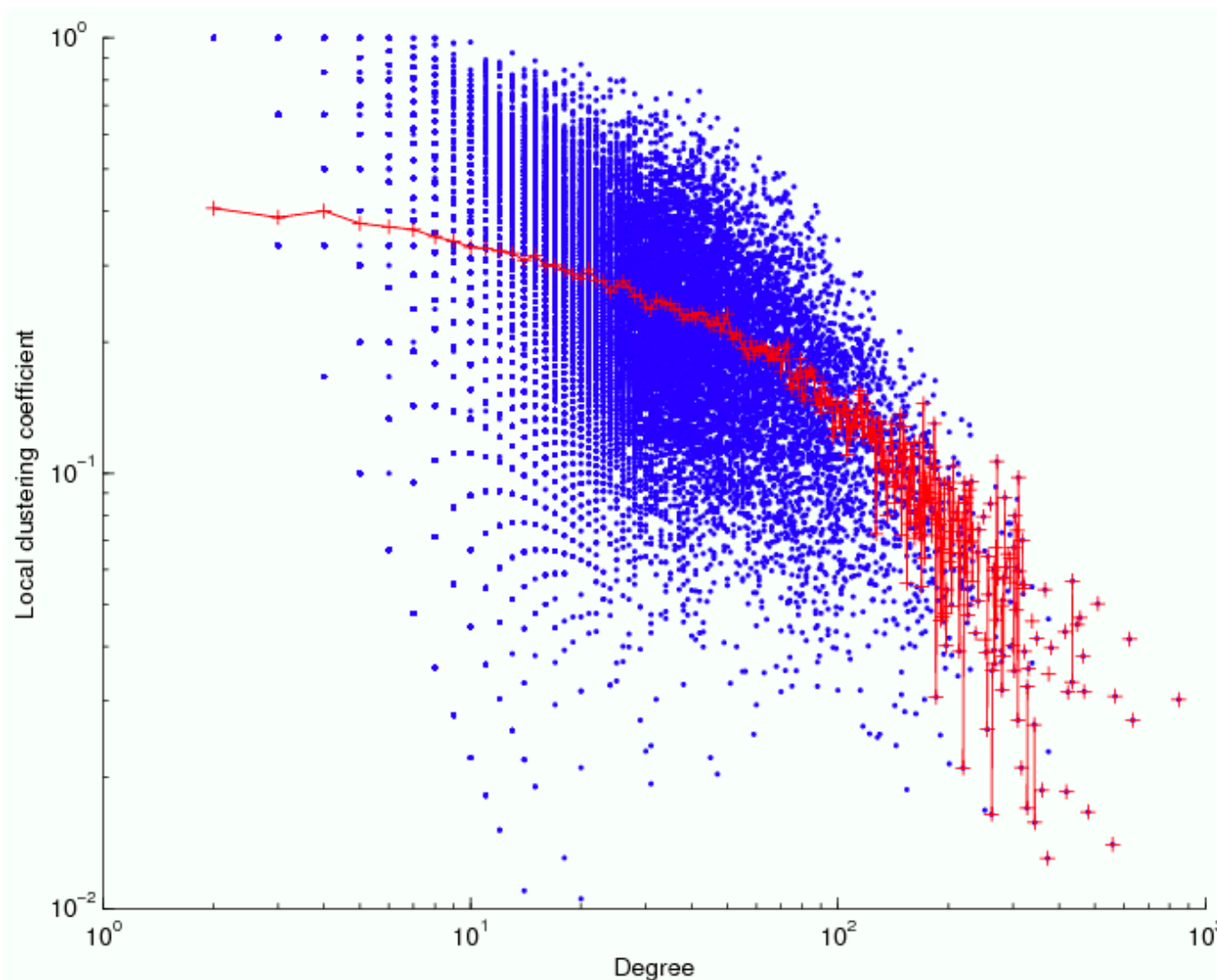
$$C_1 = 1 = \frac{6}{(4 \times 3/2)}$$



# Clustering coeff. vs degree centrality

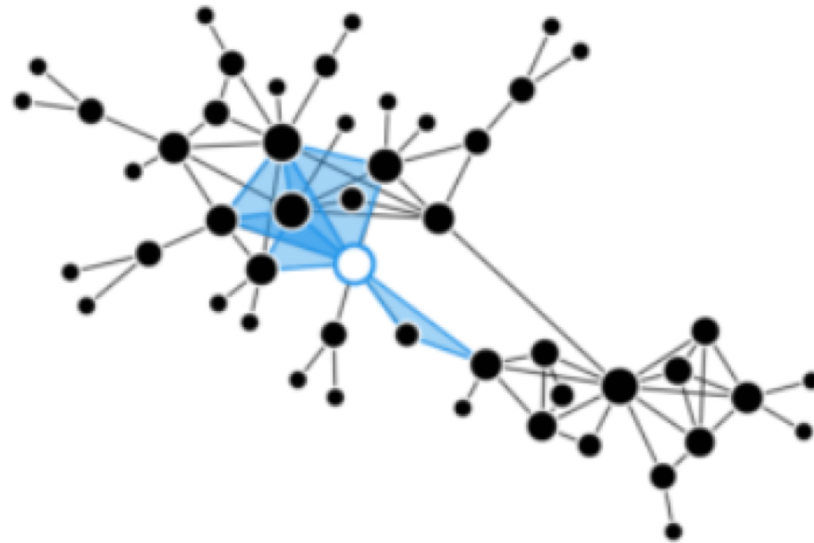
a correlation view

citation network from arXiv's High Energy Physics / Phenomenology section



when person has many friends, these friends have less edges among them, which is to be expected since a person with many friends is likely to have friends from more diverse communities, and a paper getting cited many times is likely to be cited by papers from more diverse areas

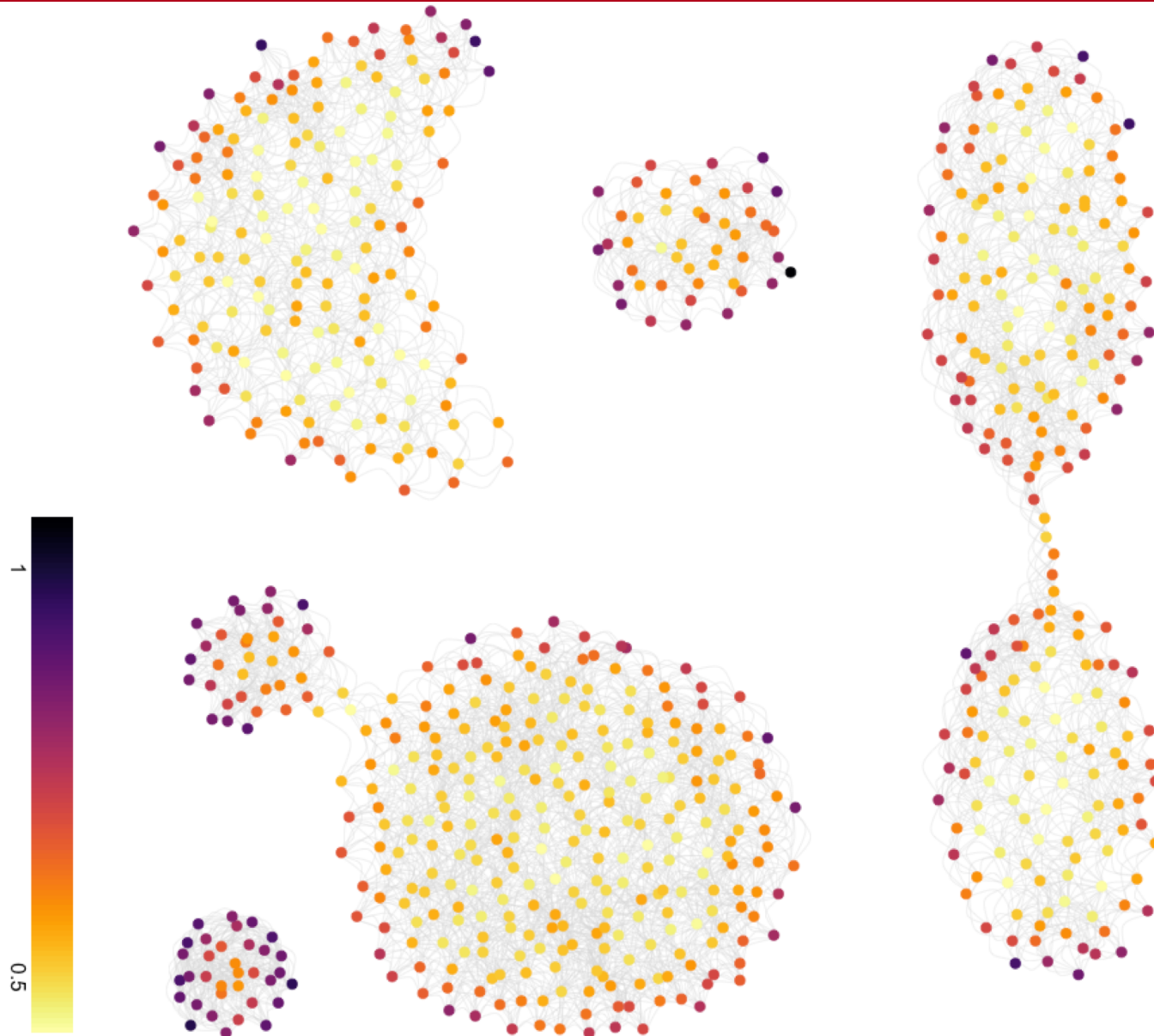




But clustering coefficient is generally hard to see and visual interpretation is considered unreliable



# Visual example





- ❑ Closeness, betweenness and clustering coefficient are **alternative** centrality measures that have a different view wrt PageRank
- ❑ They provide **useful insights** especially in social networks, as they are linked to sociology concepts
- ❑ Closeness and betweenness are based on distances, that require algorithms that are **less scalable** than PageRank
- ❑ Exploit their potential at your best

# Wrap-up

on centrality measures



Centrality measure	Technical property	Meaning
Degree (in/out)	Measures number (and quality) of <b>direct</b> connections	Cohesion <b>Entrepreneurship</b>
Attractiveness	Measures the speed of growing of a node's degree	Dinamicity Enterprising
PageRank (authorities/hubs)	Measures number (and quality) of <b>direct and indirect</b> connections	Cohesion Entrepreneurship <b>Similarity</b> /Friendship with a direction → <b>Dependence</b>
Closeness	Measures length of shortest paths	Visual centrality Significant <b>spreading points</b> Outliers/ <b>Ostracism</b>
Betweenness	Measures number of shortest paths	Brokerage <b>Structural holes</b>
Clustering coeff.	Measures number of triadic closures	Centrality in a community Cohesion of the neighbourhood



## Visual analysis

Overall organisation

Clusters (highly connected)

Sparse areas (less connected)

Cliques and strongly connected components

Disconnected components

Center/Periphery

