

# 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



#### The notion of centrality

In Network Science

## 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

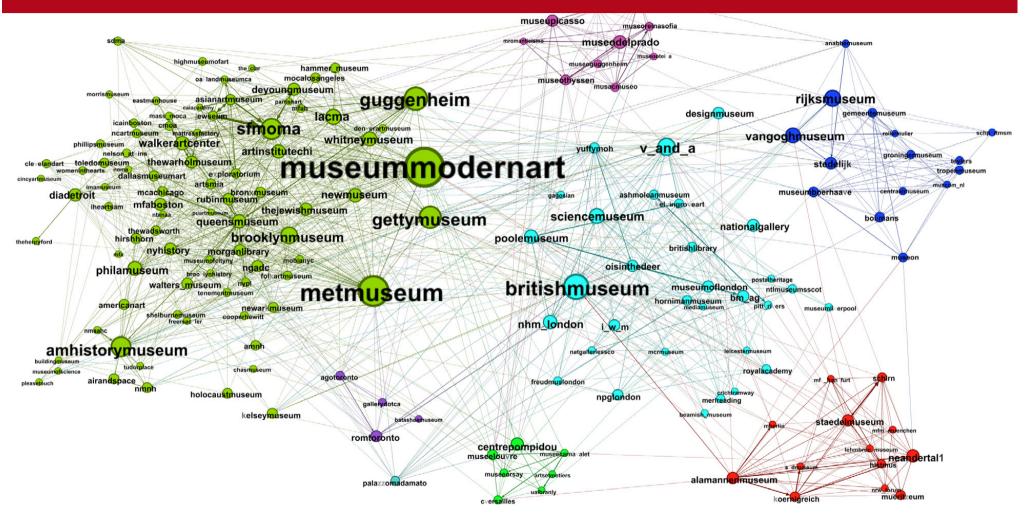
Closeness centrality [edit]

2



## An example of node centrality

museums network



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

# Degree centrality

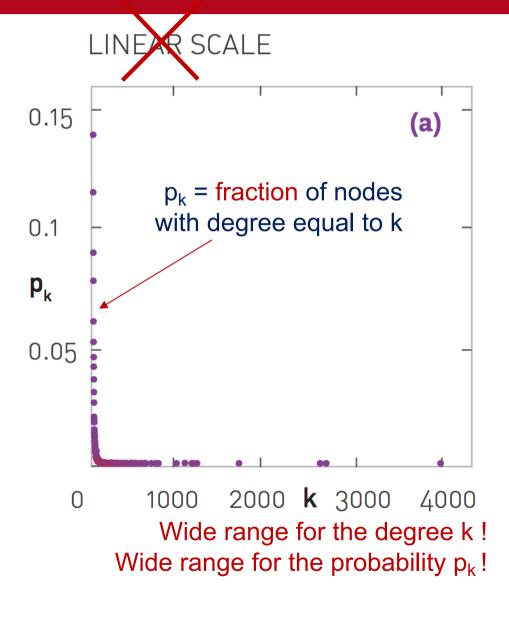
Counting the in/out degrees of nodes

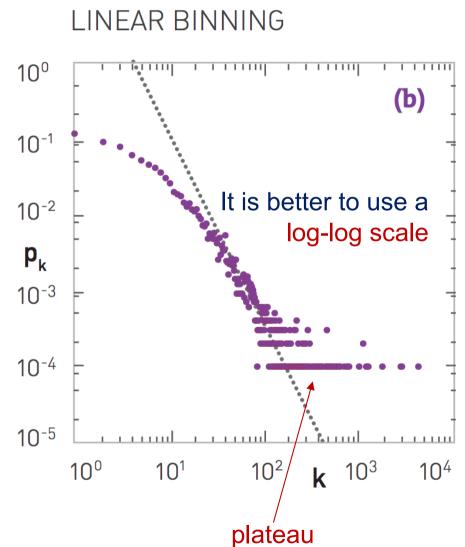


#### The degree distribution

for an undirected network

6

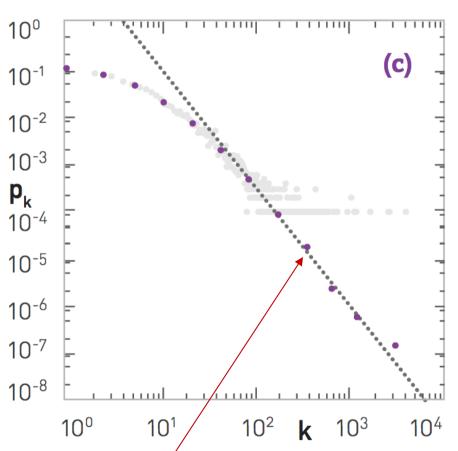




#### Alternative log representations

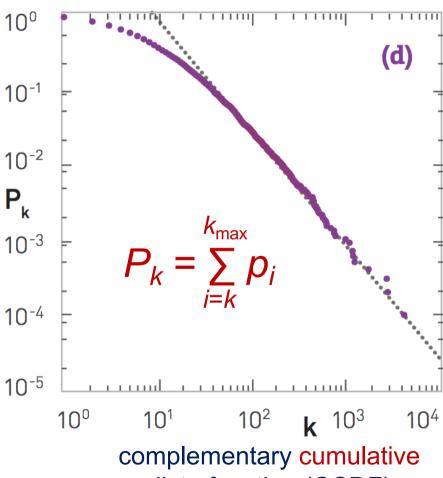
for an undirected network

#### LOG-BINNING



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

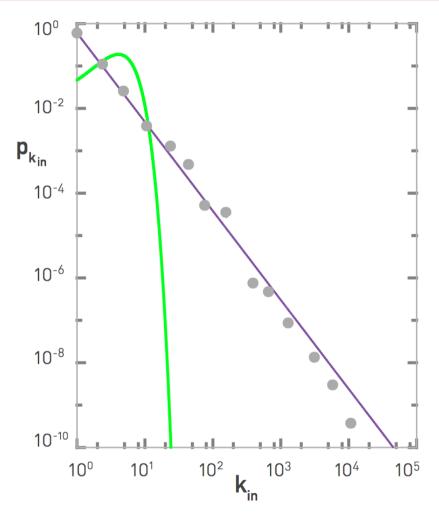
#### **CUMULATIVE**



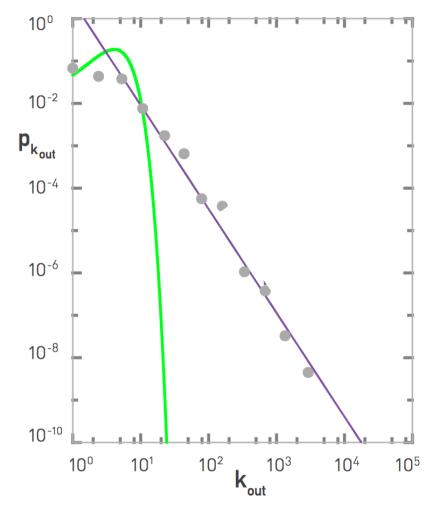


#### Two degree distributions

for directed networks



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



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

#### Pseudocode example

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

```
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
```

```
Degree Distribution

10<sup>-1</sup>

2 10<sup>-2</sup>

10<sup>-3</sup>

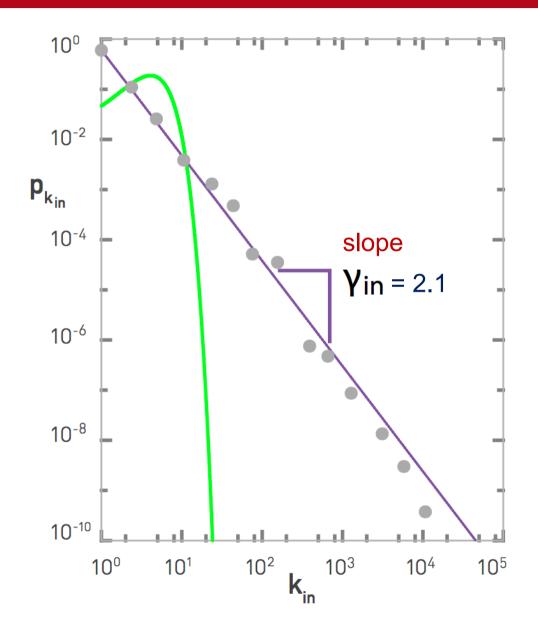
k
```

```
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

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

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

how to correctly estimate the slope y?



## 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 
$$p(k|y) = (y-1)/k_{min} \cdot (k/k_{min})^{-y}$$



# ML estimate for the exponent $\gamma$ the most reliable approach

ML criterion: find the γ that best fits the data

$$\max_{y} \sum_{i} \ln p(k_i|y)$$

where  $k_i$  is the measured degree of node i

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

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

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

#### Pseudocode example

to estimate the exponent



```
discard samples in
                               the saturation
                               region
10-3
\mathbf{p}_{\mathbf{k}}
10<sup>-5</sup>
            choose an
           appropriate
10^{-7}
              k_{min} = 49
10<sup>-9</sup>
        10°
                                       10^{2}
                                                      10<sup>3</sup>
                                                                    104
                       10<sup>1</sup>
```

```
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 $\gamma$ ∈ [2,5]

192,244	609,066	6.34			
205 720		0.04	-	-	3.42*
325,729	1,497,134	4.60	2.00	2.31	_
4,941	6,594	2.67	_	-	Exp.
36,595	91,826	2.51	4.69*	5.01*	_
57,194	103,731	1.81	3.43*	2.03*	_
23,133	93,439	8.08	_	-	3.35*
702,388	29,397,908	83.71	_	_	2.12*
449,673	4,689,479	10.43	3.03**	4.00*	_
1,039	5,802	5.58	2.43*	2.9 0*	_
2,018	2,930	2.90	_	_	2.89*
	36,595 57,194 23,133 702,388 449,673 1,039	36,595       91,826         57,194       103,731         23,133       93,439         702,388       29,397,908         449,673       4,689,479         1,039       5,802	36,595       91,826       2.51         57,194       103,731       1.81         23,133       93,439       8.08         702,388       29,397,908       83.71         449,673       4,689,479       10.43         1,039       5,802       5.58	36,595       91,826       2.51       4.69*         57,194       103,731       1.81       3.43*         23,133       93,439       8.08       -         702,388       29,397,908       83.71       -         449,673       4,689,479       10.43       3.03**         1,039       5,802       5.58       2.43*	36,595       91,826       2.51       4.69*       5.01*         57,194       103,731       1.81       3.43*       2.03*         23,133       93,439       8.08       -       -         702,388       29,397,908       83.71       -       -         449,673       4,689,479       10.43       3.03**       4.00*         1,039       5,802       5.58       2.43*       2.9 0*

<sup>\* =</sup> good statistical fit with a power-law

<sup>\*\* =</sup> good fit for a power-law with an exponential cutoff

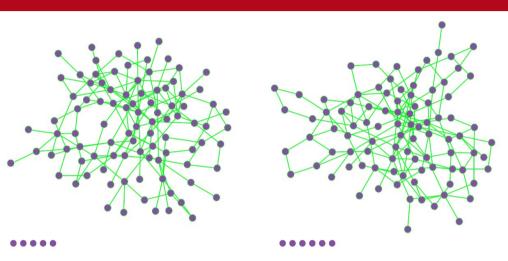
Exp = good fit with an exponential distribution e<sup>-ak</sup>

# Explaining the power-law

Preferential attachment



#### Random networks Erdös-Rényi model 1959/60





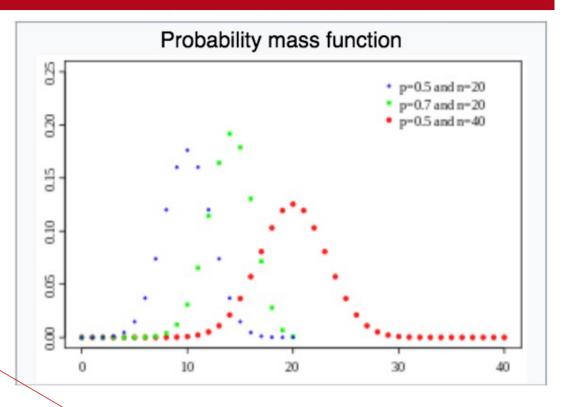
- ☐ The random network is the simplest model:
  - pick a probability p, with 0 activate each link <math>(i,j) with probability p
- ☐ The number of links is variable
- ☐ There might be isolates
- Easy to calculate fundamental parameters



#### Binomial distribution

explains the degree distribution for random networks

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



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

## Degree distribution

in random networks

☐ 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

this defines p

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

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

**Variance** 

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

$$\sigma_{x}^{2} = (N-1)p(1-p) \simeq \langle k \rangle$$

tight around the mean

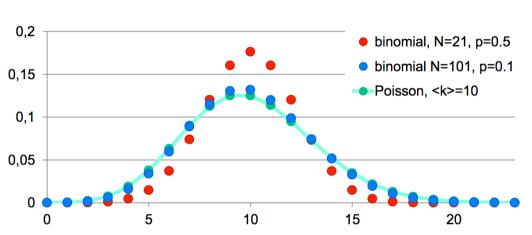


#### Poisson approximation

why random networks are called Poisson networks

Poisson distribution (easier to use)

$$P\left[x=k\right] = \frac{m_x^k}{k!} \cdot e^{-m_x} \quad _{\text{0,2}}$$



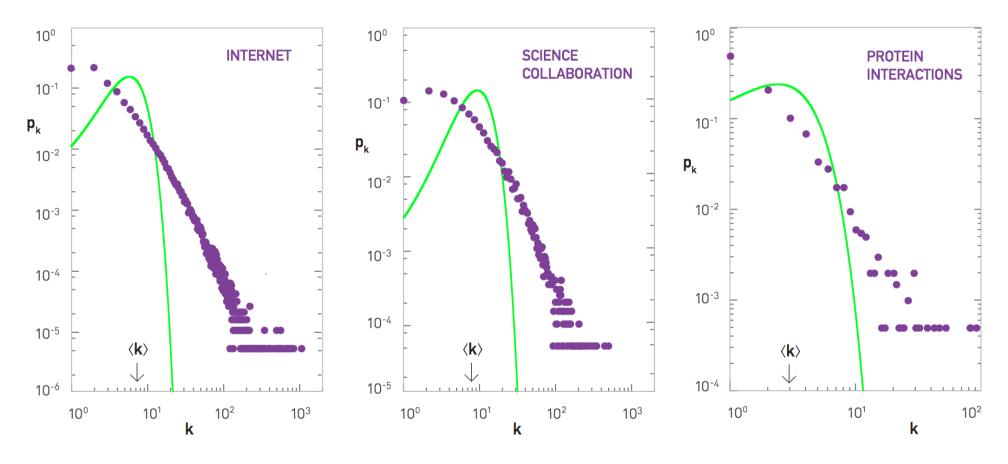
☐ Very good approximation of binomial for small p (and at small k) active part p

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



#### Are real networks Poisson?

no, they aren't



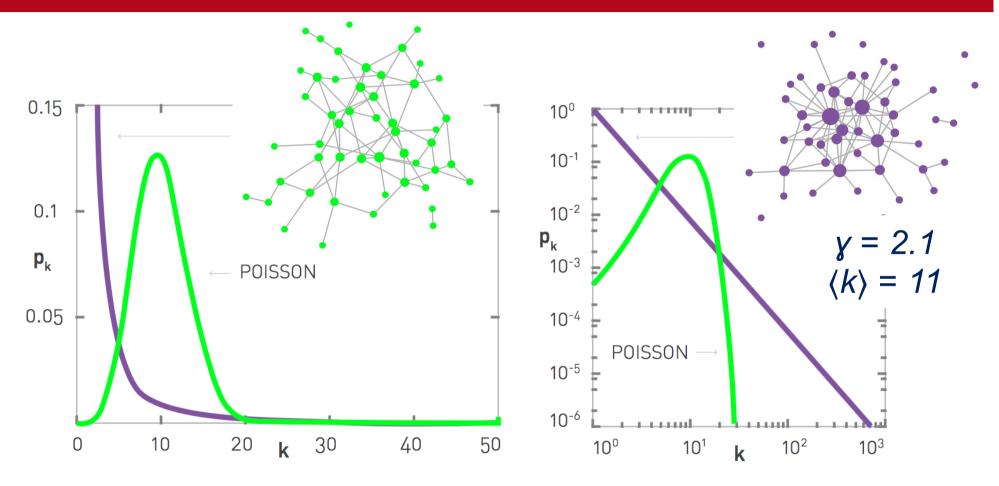
No! Poisson networks are deprived of hubs

... 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



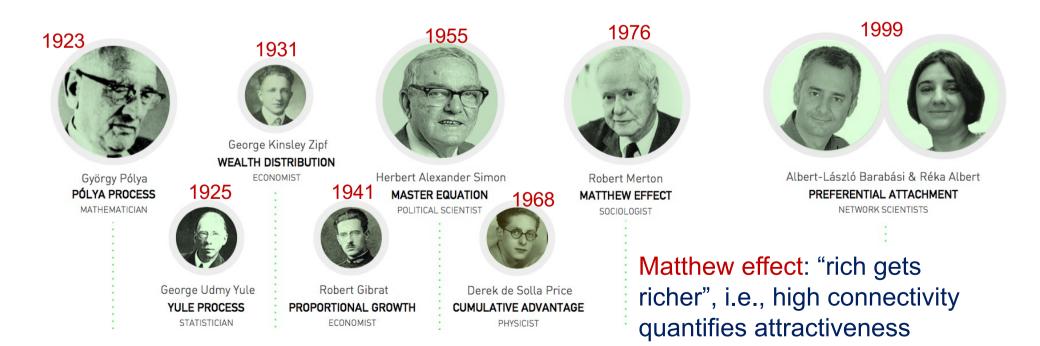
#### Preferential attachment

a simple concept that (partially) explains the power-law

#### Nodes link to the more connected nodes

e.g., think of www

#### This idea has a long history



#### The Barabasi-Albert model

Barabási, Albert. "Emergence of scaling in random networks" (1999)

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

□ Growth

add a node (the Nth) with m links that connect the node to nodes in the network

Preferential attachment

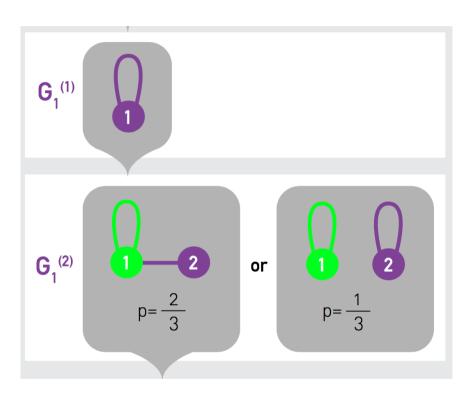
 $p_i = k_i/C$  probability of connecting to node i

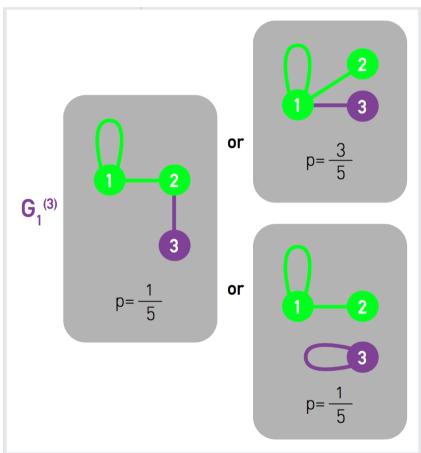
$$p_i = 1/C$$
 for self-loops

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



# An example with m=1







## Approximate analysis

evolution of nodes degree

☐ Increase in the degree (at each step)

$$\Delta k_i \simeq m \cdot k_i / (1+2m(N-1)) \simeq k_i / 2N$$
trials probability per trial

Approximation in the continuous domain

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

Integration

$$ln(k_i) = \frac{1}{2} ln(N) + cost. \rightarrow k_i = c N^{\frac{1}{2}}$$

 $\square$  Recalling that node *i* joins the network at time N = i

$$k_i(N=i) = m \rightarrow k_i(N) = m (N/i)^{\frac{1}{2}}$$
<sup>1/2</sup> is the dynamic exponent

# Approximate analysis degree distribution

- $\Box$  The number of nodes with degree smaller than k is

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

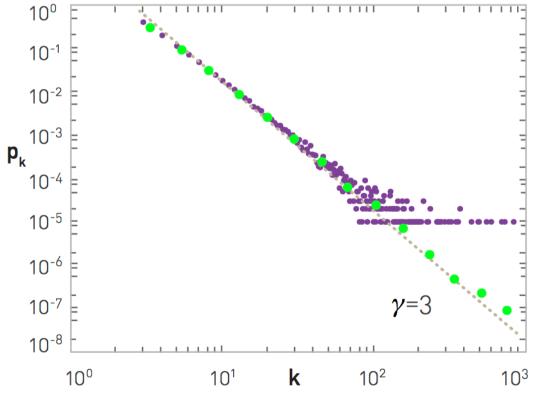
- □ CDF is  $P_k = P[k_i \le k] = 1 (m/k)^2$
- The degree distribution is

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

## The Barabasi-Albert model

wrap up

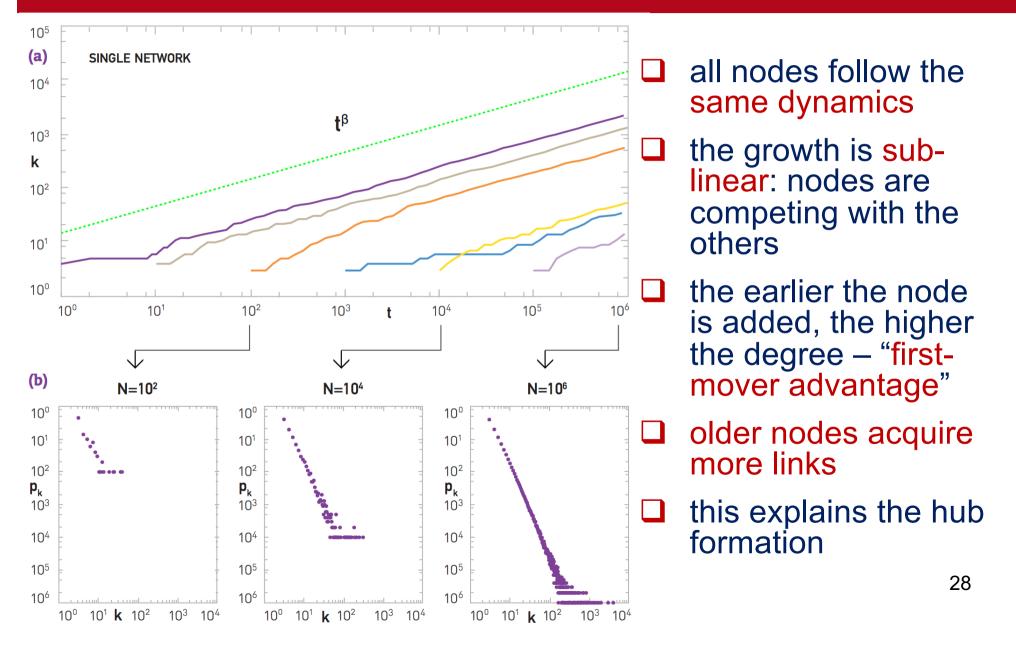
- Depending on the implementation there might be self/multiple links
- Most nodes have a small degree (exactly m for the youngest ones)
- Hubs appear
- The average degree is  $\langle k \rangle = 2m$ , and in fact  $L = Nm = \frac{1}{2} \langle k \rangle N$
- The resulting degree distribution is always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





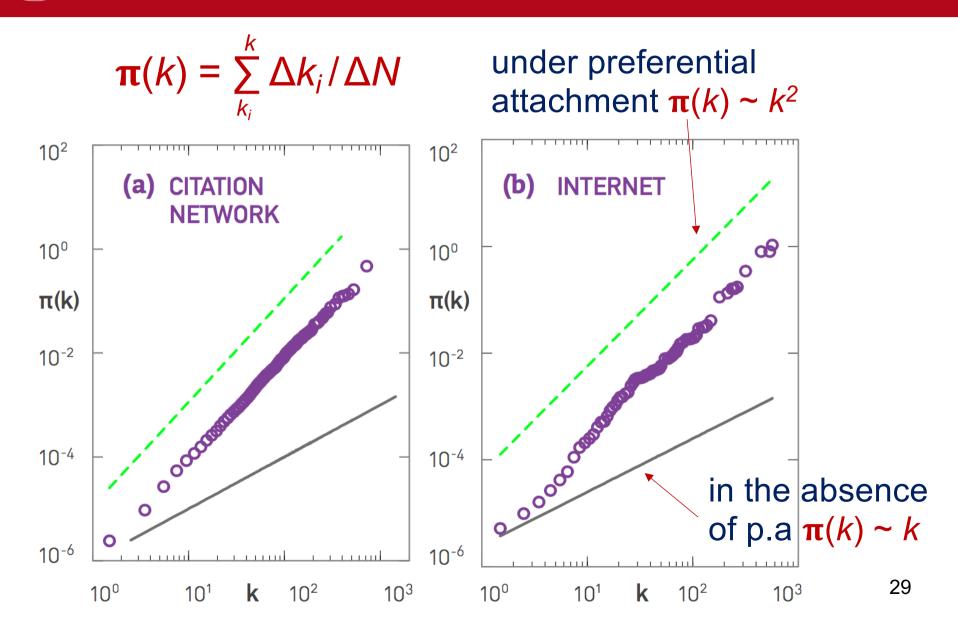
#### The Barabasi-Albert model

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



## Measuring preferential attachment

in real networks





#### The Bianconi-Barabasi model

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

#### The model:

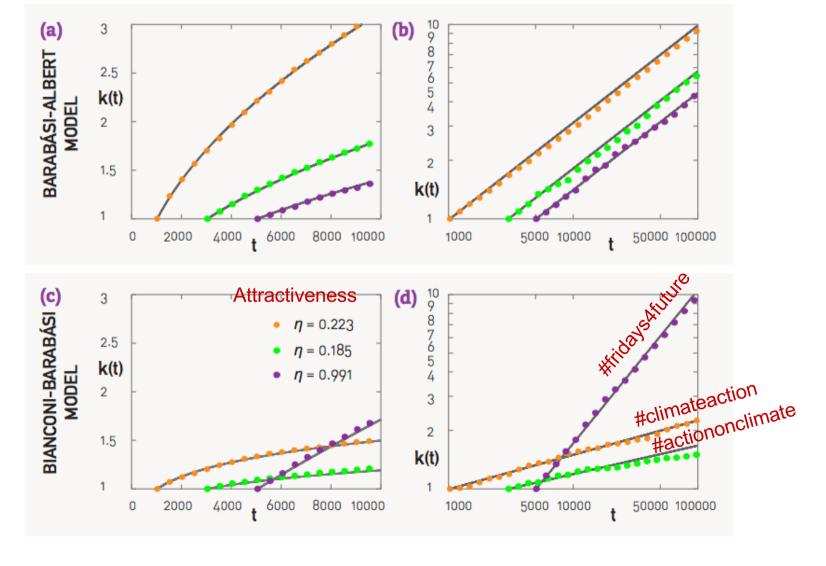
- □ Growth at time step N a new node i=N is added with m links and fitness  $\eta_i$
- Attractiveness (or fitness) is a random number drawn from a given distribution  $\rho(\eta)$  a quality of the individual to attract links
- Preferential attachment probability of linking to node i is proportional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



#### An example

properties of the Bianconi-Barabasi model

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





#### Approximate analysis starting point

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

trials probability per trial

- □ Increase in the degree  $\Delta k_i \simeq m \cdot k_i \eta_i / \sum k_i \eta_i$
- lacksquare We show that  $\sum k_i \eta_i \simeq m \ N \cdot C$  (see proof)



## Approximate analysis

the denominator

- $\square$  Analysis of denominator  $\sum k_i \eta_i$ 
  - $\rightarrow$  average value wrt  $\eta$
  - $\rightarrow$  hypothesis  $k_i \simeq m (N/i)^{\beta(\eta i)}$

- Swap integrals  $A \simeq \int m N^{\beta(\eta)} \left[ \int_{1}^{\eta} i^{-\beta(\eta)} di \right] \eta \cdot \rho(\eta) d\eta$
- Integrate constant C  $A \simeq m \ N \cdot \int (1 N^{\beta(\eta)-1}) \ \eta \ \rho(\eta) \ d\eta$   $1-\beta(\eta)$



### Approximate analysis

evolution of nodes degrees

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

- Increase in the degree  $\Delta k_i \simeq m \cdot k_i \eta_i / \sum k_i \eta_i$
- It is  $\sum k_i \eta_i \simeq m N \cdot C$

#### Hence:

By inspection of the above

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

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

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

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

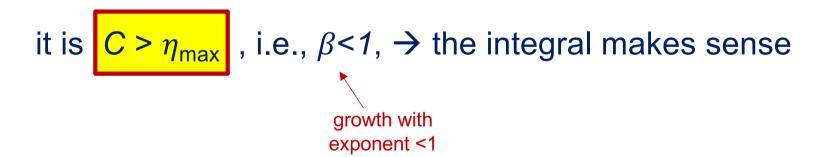


# Approximate analysis constant C

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

$$C = \int \frac{\eta \, \rho(\eta) \, d\eta}{1 - \beta(\eta)} \rightarrow 1 = \int_{0}^{\eta_{\text{max}}} (C - \eta)^{-1} \, \eta \, \rho(\eta) \, d\eta$$

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



it also is  $C \le 2\eta_{\text{max}}$ 

# Approximate analysis degree distribution

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

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

#### The degree distribution is

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

# Equal fitness the Barabasi-Albert model

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

Coefficient C = 2 since  $\int_{0}^{\eta_{\text{max}}} (C/\eta - 1)^{-1} \delta(\eta - 1) d\eta = (C - 1)^{-1} = 1$ 

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

Degree distribution

$$p_{k} = C \int_{0}^{\eta_{\text{max}}} \eta^{-1} m^{C/\eta} k^{-(C/\eta+1)} \delta(\eta-1) d\eta = 2 m^{2} k^{-3}$$



# Uniform fitness the model

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

- □ Coefficient C = 1.255 since  $\int_{0}^{1} (C/\eta 1)^{-1} d\eta = 1 \implies e^{-2/C} = 1-1/C$
- Exponential degree  $k_i \simeq m (N/i)^{\eta i/C}$
- ☐ Each node has its own dynamic exponent !!!

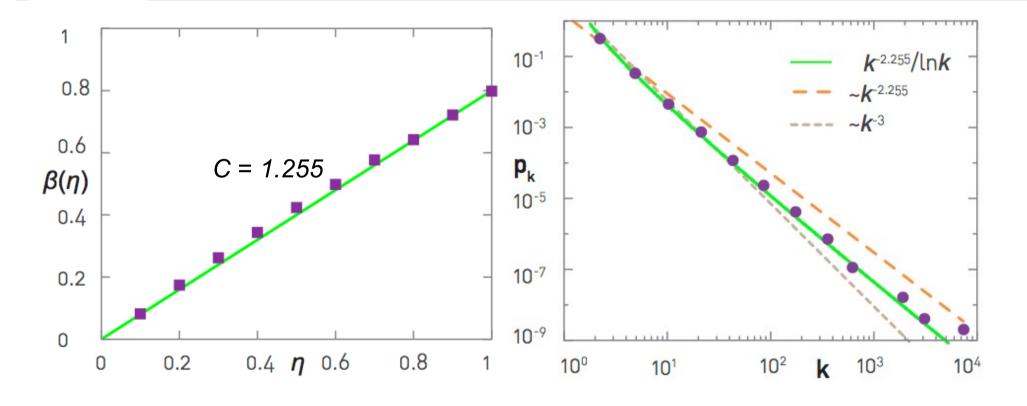
#### Degree distribution

$$p_{k} = C/k \int_{0}^{1} \eta^{-1} e^{-C \ln(k/m)/\eta} d\eta \sim k^{-(1+C)} / \ln(k)$$

$$e^{-b} - b E_{1}(b), b = C \ln(k/m)$$
exponential integral E<sub>1</sub>

#### Uniform fitness

the measured data



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

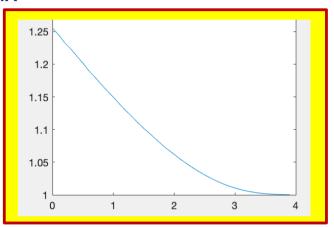


#### Exponential fitness the model

What if  $\rho(\eta)$  = a e<sup>-a $\eta$ </sup> / (1-e<sup>-a</sup>) and  $\eta_{\text{max}}$ 

□ C rapidly converges to C=1  $\int_{0}^{1} (C/\eta - 1)^{-1} \rho(\eta) d\eta = 1$ 

$$\int_{0}^{1} (C/\eta - 1)^{-1} \rho(\eta) d\eta = 1$$



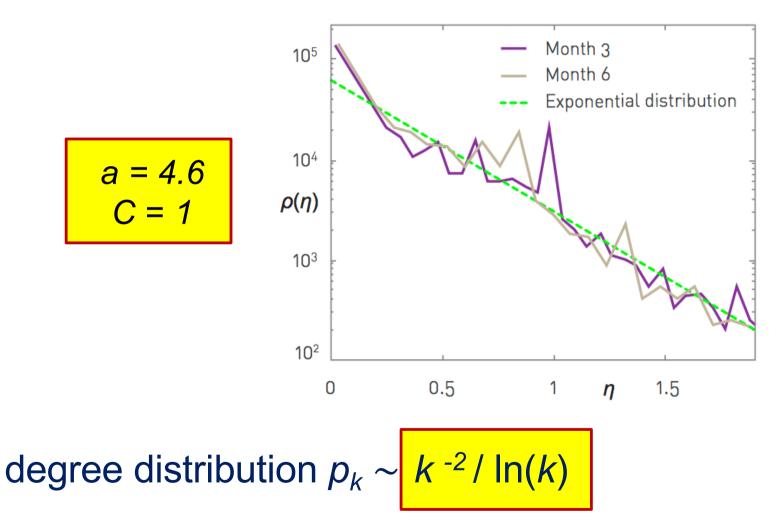
- Exponential degree  $k_i \simeq m (N/i)^{\eta i/C}$
- Each node has its own dynamic exponent !!!

Degree distribution

$$p_{k} = C/k \int_{0}^{1} \eta^{-1} e^{-C \ln(k/m)/\eta} \rho(\eta) d\eta \sim \frac{k^{-(1+C)}/\ln(k)}{\text{exponential integral E}_{1}}$$

# Exponential fitness

the www



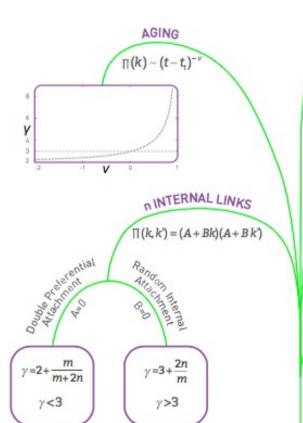


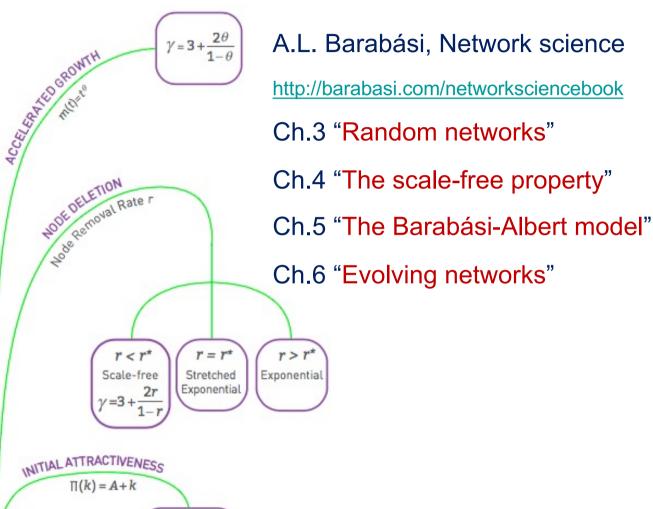
#### Other ideas for extension

of the Albert-Barabasi model

#### **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.





p. - (k+A)-

# Properties of the power-law

scale-free and random networks



# The largest hub

natural cutoff under the power-law

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

The size of the largest hub is captured by

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

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

is the natural cutoff it explains large hubs

<k>

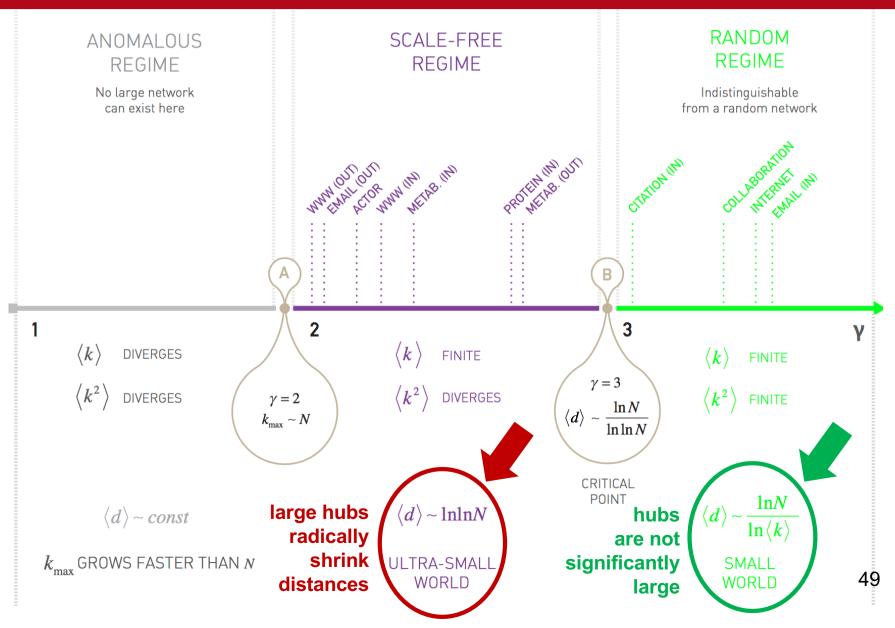


They diverge with N if  $\gamma < n+1$ mean (n=1) doesn't diverge for  $\gamma \ge 2$ variance (n=2) diverges for  $\gamma < 3$ and the network does not have a scale

(scale-free regime)



# The scale-free regime



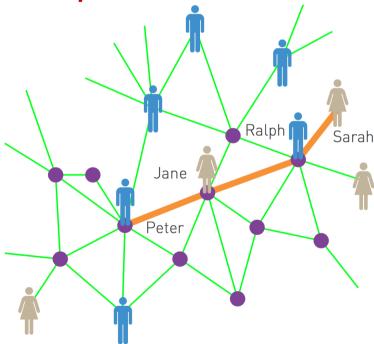


## Small world property

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

# Distances in random graphs

theoretical result

- 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

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

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

e.g.: on earth we are  $N=7\cdot10^9$  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

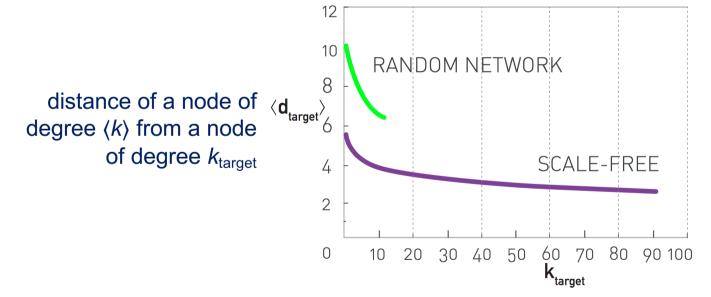


#### Distances in scale-free networks

the ultra-small-world

■ The average distance increases as ln(ln(N)), much slower than N or ln(N)

e.g. in www  $N=7.10^9$ , ln(N)=22.7, ln(ln(N))=3.12 (very small)



□ The large hubs radically shrink the distance between nodes → 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
- $\square$  # of friends (in the average) =  $\langle k \rangle$
- ☐ # of friends of friends ~ 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 γ, but please use the ML approach: γ 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?

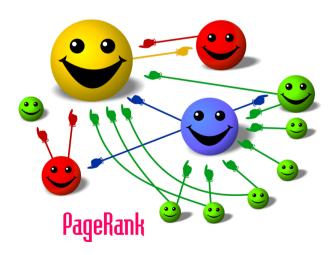
links as votes

- □ 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



# The Google's view quoting Google

- □ 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

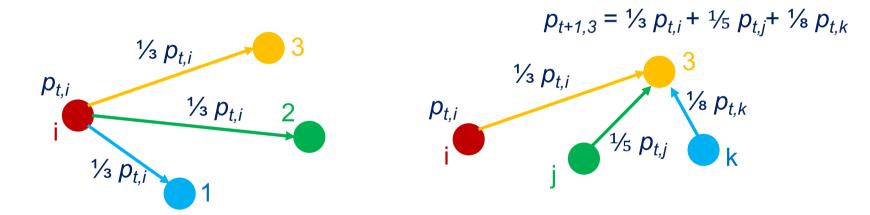




#### A random walk on www

the rationale behind PageRank

- $\Box$  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

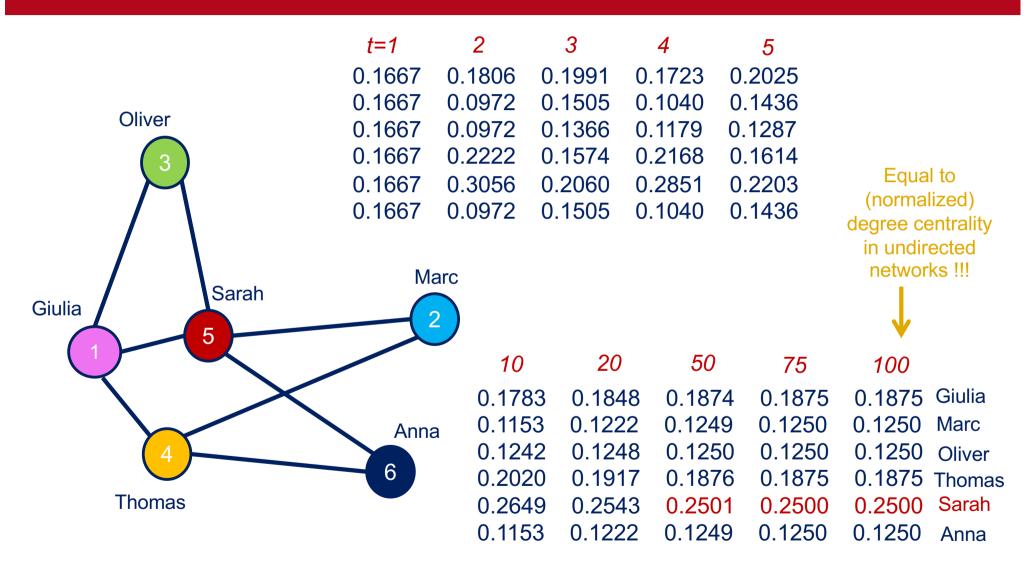


- this identifies a Markov chain
- after a while probabilities settle to a steady state = the PageRank vector



### Example

of the random walk effect on a friends' network



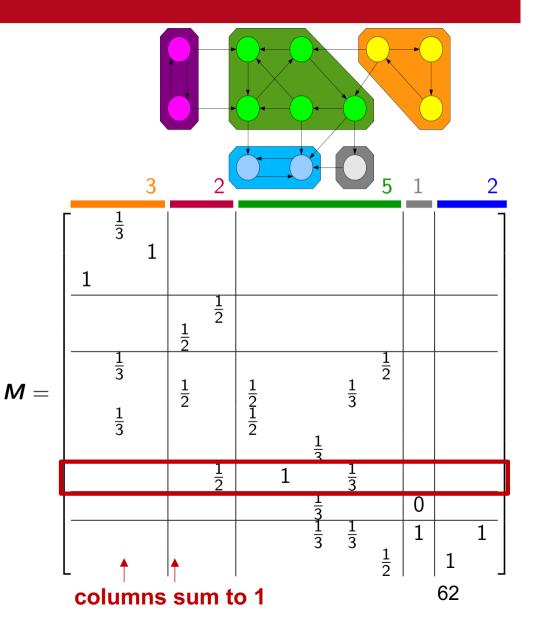


#### Matrix formalization

of the random walk

- p<sub>t</sub> stochastic vector (positive entries which sum up to 1)
- M normalized adjacency matrix (column stochastic)

- $p_{\infty} = M p_{\infty}$  converges to an eigenvector of M (with eigenvalue 1)
- $\mathbf{p}_{\infty} = \mathbf{d}$  for undirected networks where  $\mathbf{A} = \mathbf{A}^{\mathsf{T}}$



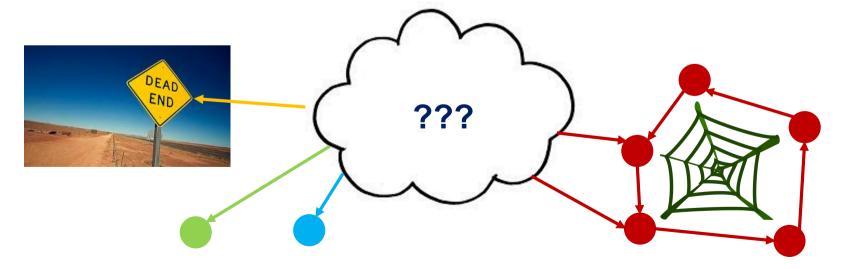


# Problems in the random walk

dead ends and spider traps

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



#### **Teleportation**

as a method to overcome problems

#### 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

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  $\mathbf{q}$  independent of the node she/he is in, e.g.,  $\mathbf{q} = 1/N$ for uniform probability



# PageRank with restart or simply PageRank

dead ends

no dead ends

normalization

no spider traps

Markov chain

PageRank equation

$$\mathbf{A} = \mathbf{A}_0 + \mathbf{b} \mathbf{e}^{\mathsf{T}}$$
 indicating vector of dead ends

$$M = A \operatorname{diag}^{-1}(d), \qquad d = A^{T}1$$

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

$$\mathbf{p}_{t+1} = \mathbf{M}_1 \mathbf{p}_t$$

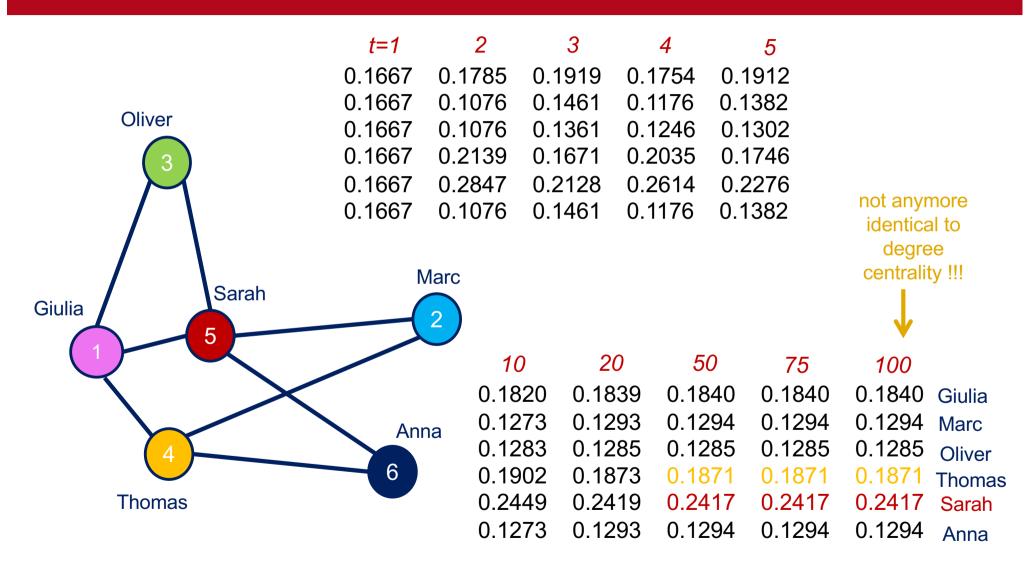
PageRank centrality vector

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



### Example

#### of PageRank with restart on a friends' network

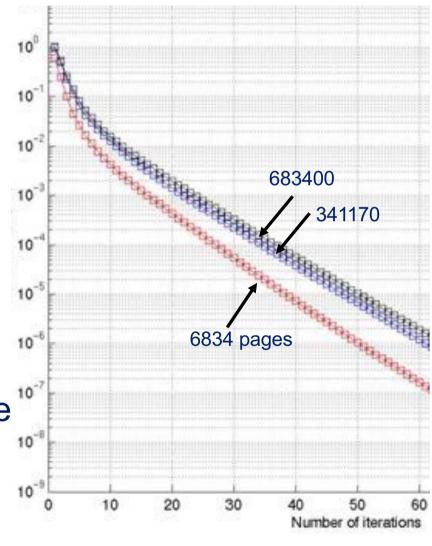




# Convergence properies of PageRank

an overview

- It corresponds to the stationary behaviour of the Markov chain
- $\square$   $p_{\infty}$  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  $\mathbf{c}$ )
- $oldsymbol{\square}$   $oldsymbol{p}_{\infty}$  converges in few iterations, typically  $oldsymbol{p}_{40} \simeq oldsymbol{p}_{\infty}$





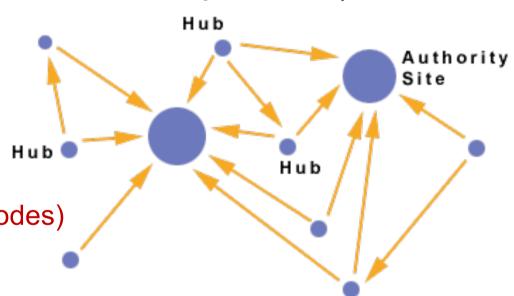
# Hubs and Authorities

what we can get from PageRank

## □ 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  $\mathbf{A}_0^{\mathsf{T}}$  (related to the out-degree of nodes)

#### authority or hub?



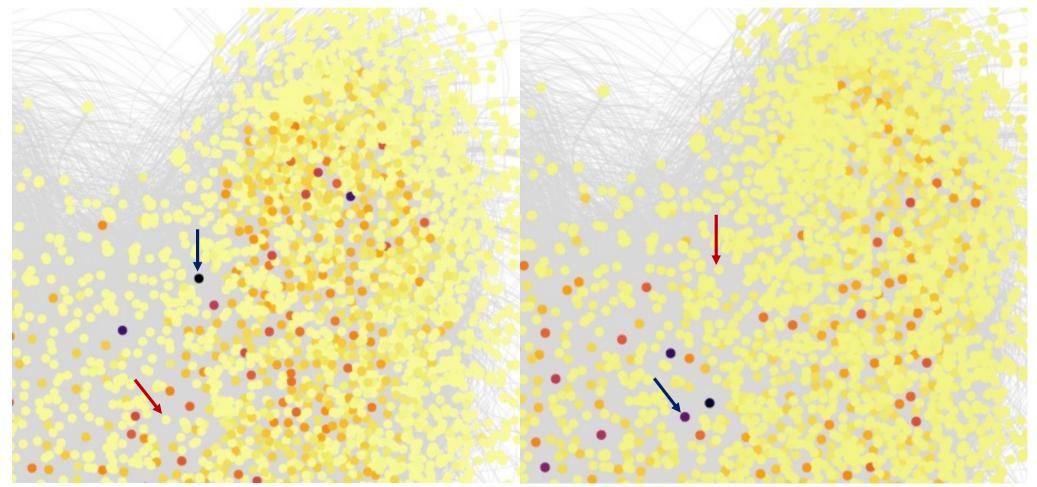


# Example of PageRank centrality

wikipedia administrator elections and vote history data <a href="https://snap.stanford.edu/data/wiki-Vote.html">https://snap.stanford.edu/data/wiki-Vote.html</a>

Hubs

Authorities





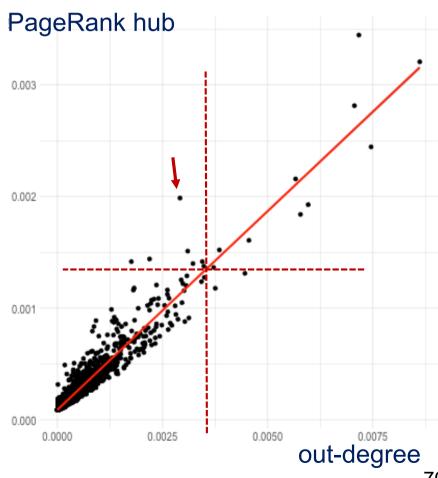
### PageRank versus degree centrality

wikipedia administrator elections and vote history data

#### **Authorities**

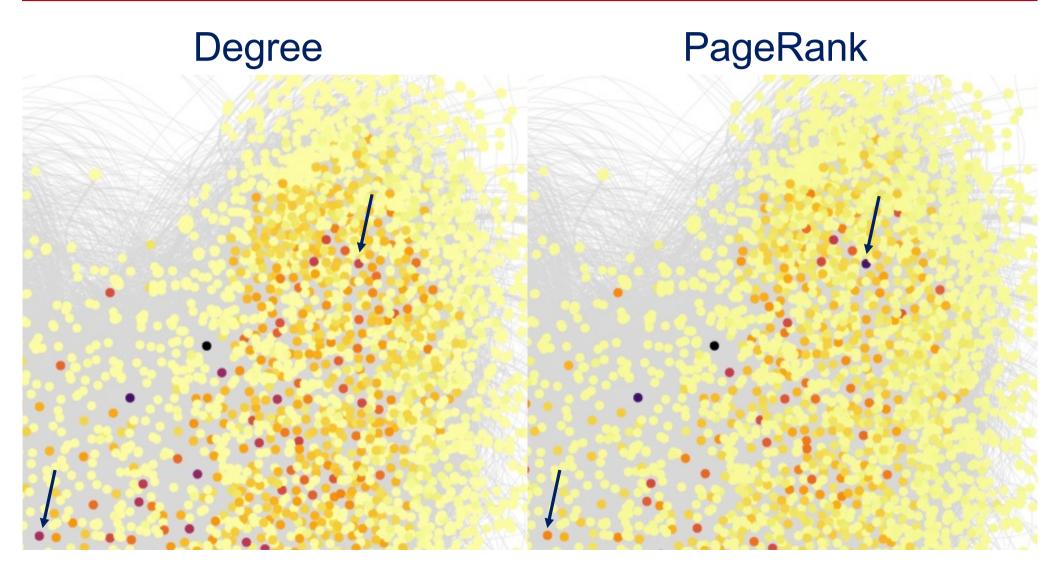
# PageRank authority 0.004 0.003 0.002 0.001 0.003 0.004 0.002 0.001 in-degree

#### Hubs



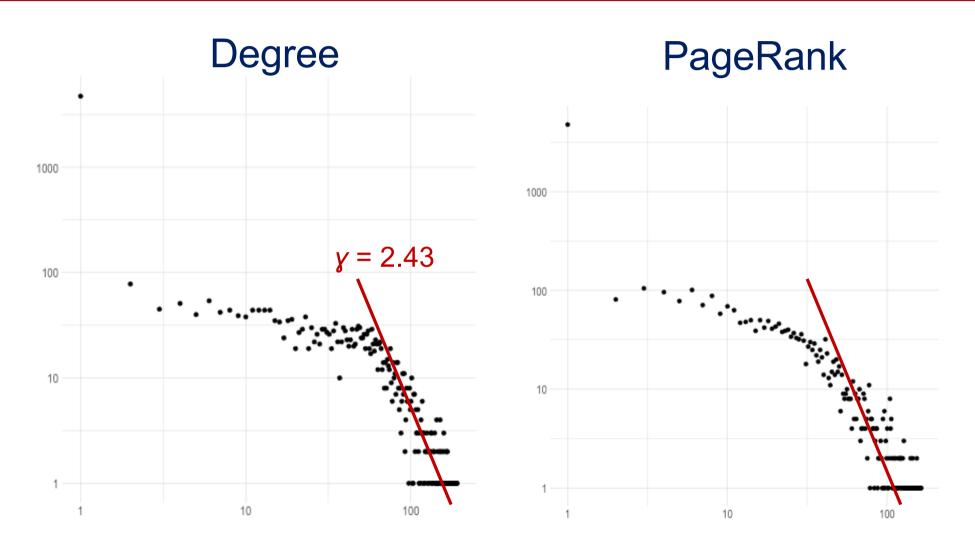


# PageRank versus degree authorities



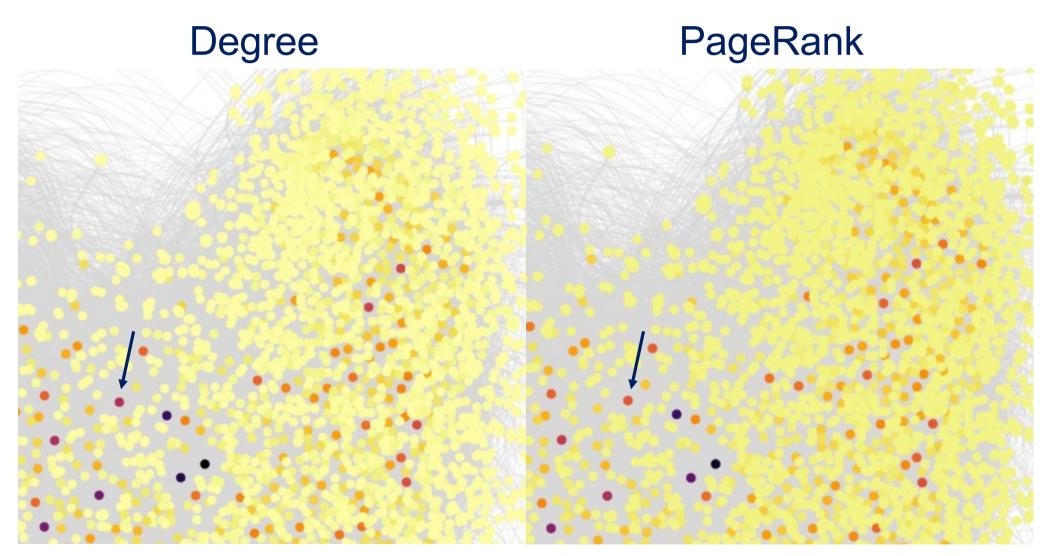


# PageRank versus degree authorities



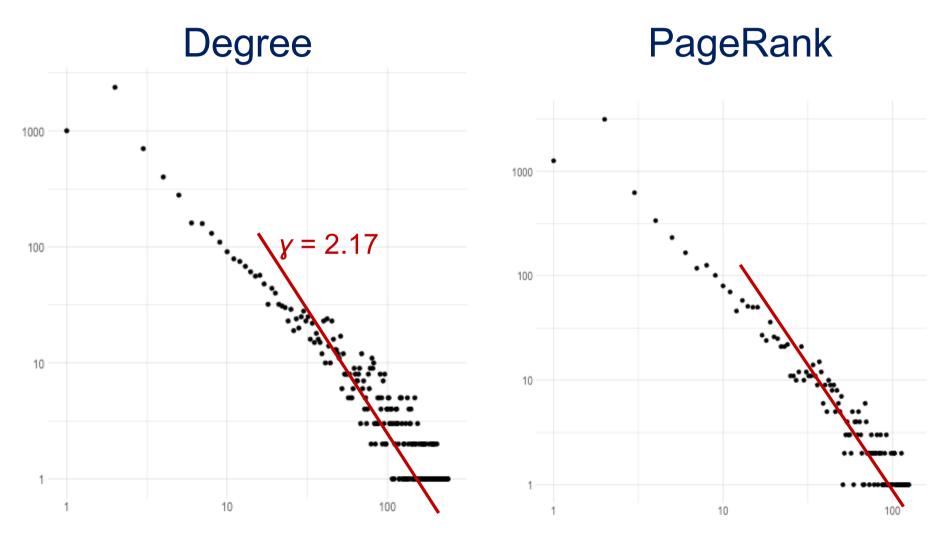


# PageRank versus degree hubs





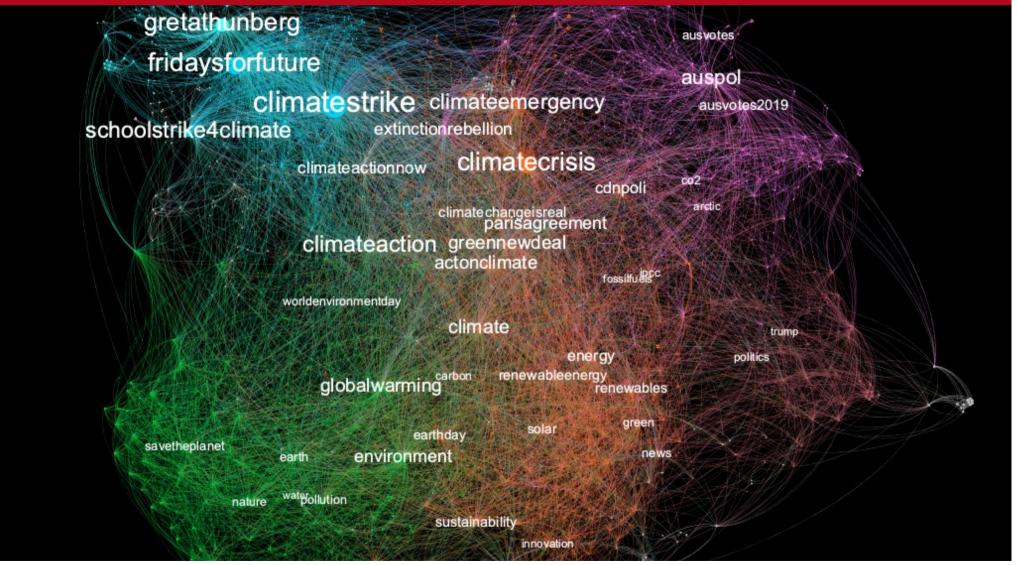
# PageRank versus degree hubs





### 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



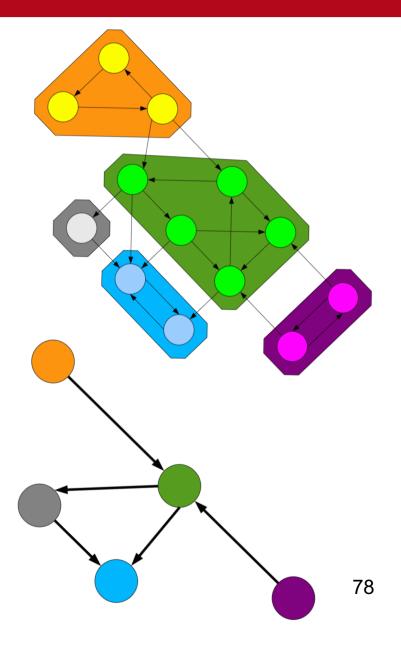
# Convergence properties

of PageRank power iterations

### The condensation graph

ordering an adjacency matrix

- Strong connectivity induces a partition in disjoint strongly connected sets  $V_1, V_2, ..., V_K$
- By reinterpreting the sets as nodes we obtain a condensation graph  $g^*$  where  $i \rightarrow j$  is an edge if a connection exists between sets  $\mathcal{V}_i \rightarrow \mathcal{V}_i$





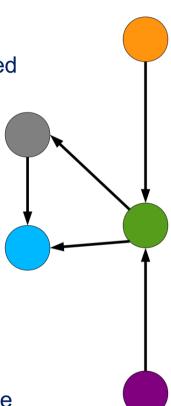
### Properties of the condensation graph

ordering an adjacency matrix

- $\Box$   $g^*$  does not contain cycles
  - otherwise the sets in the cycle would be strongly connected
- $\Box$   $G^*$  has at least one root and one leaf
  - and every node in the graph can be reached from one of the roots
- $\Box$   $g^*$  allows a particular reordering

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

procedure: identify a root n<sub>1</sub> and remove it from the network, then identify a new root; cycle until all nodes have been selected

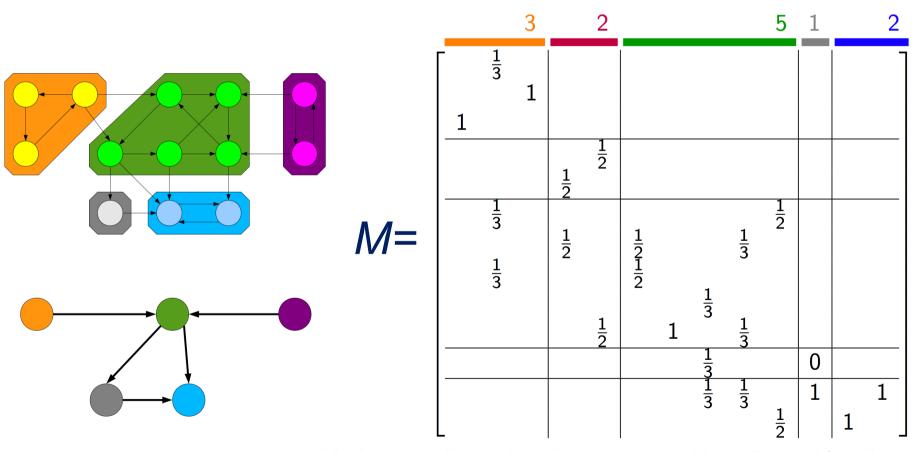




### Matrix representation

of the condensation graph

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



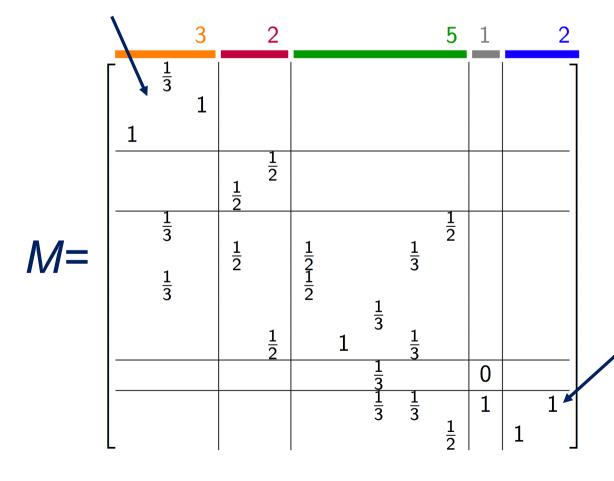
blocks in the diagonal are irreducible = no block-diagonal form! 80

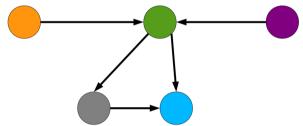


#### Perron-Frobenius theorem

of the condensation graph

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



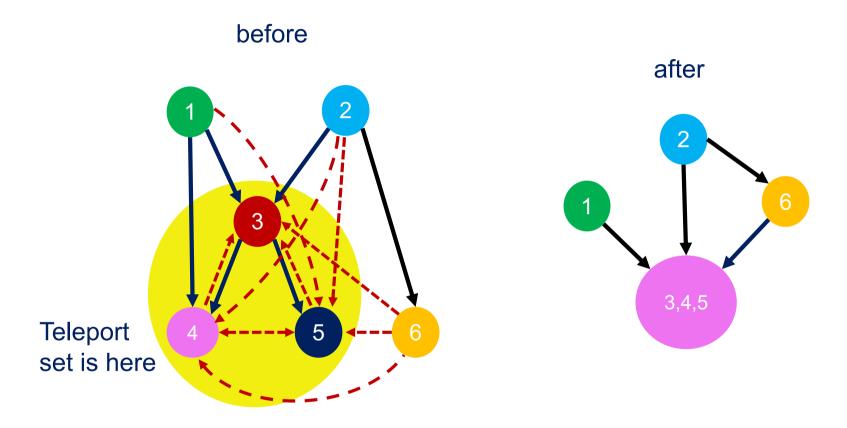


each leaf-block has at least one eigenvalue in the unit circle;  $\lambda$ =1 is always available, the others are distinct



#### The teleportation effect

it implies only one leaf



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



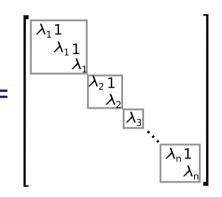
#### Lemma

on generalized eigenvectors

- □ PageRank matrix  $M_1 = c M + (1-c) q 1^T$
- □ Normalization property  $\mathbf{1}^{\mathsf{T}} \mathbf{M}_1 = \mathbf{1}^{\mathsf{T}}$
- ☐ Jordan form  $M_1 = V J V^{-1}$

carries the right (generalized) eigenvectors **e**<sub>i</sub> of **M**<sub>1</sub>

carries the eigenvalues of **M**<sub>1</sub>



$$1^{T} M_{1} V = 1^{T} V$$

$$= 1^{T} V J$$

$$\rho \text{ only one value is 0}$$

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



#### Main result

for the eigenstructure of the PageRank matrix

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



- $\square$   $M_1$  has one eigenvalue equal to 1
- □ The remaining eigenvalues satisfy  $|\lambda| \le c$

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



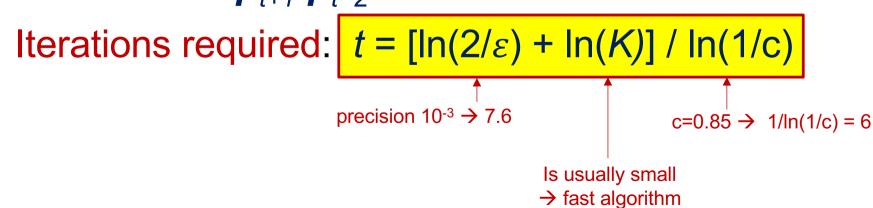
#### Convergence properties

of the PageRank power iteration

$$p_t = M_1 p_{t-1} = M_1^t p_0 = V J^t V^{-1} p_0$$

gets large for high multiplicity max eigenvalue multiplicity

- Triangular inequality:  $\|\boldsymbol{p}_{t+1}-\boldsymbol{p}_t\|_2 \lesssim 2K c^t$
- □ Precision  $\varepsilon$ :  $\|\boldsymbol{p}_{t+1} \boldsymbol{p}_t\|_2 < \varepsilon$



# 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 q = δ<sub>i</sub>

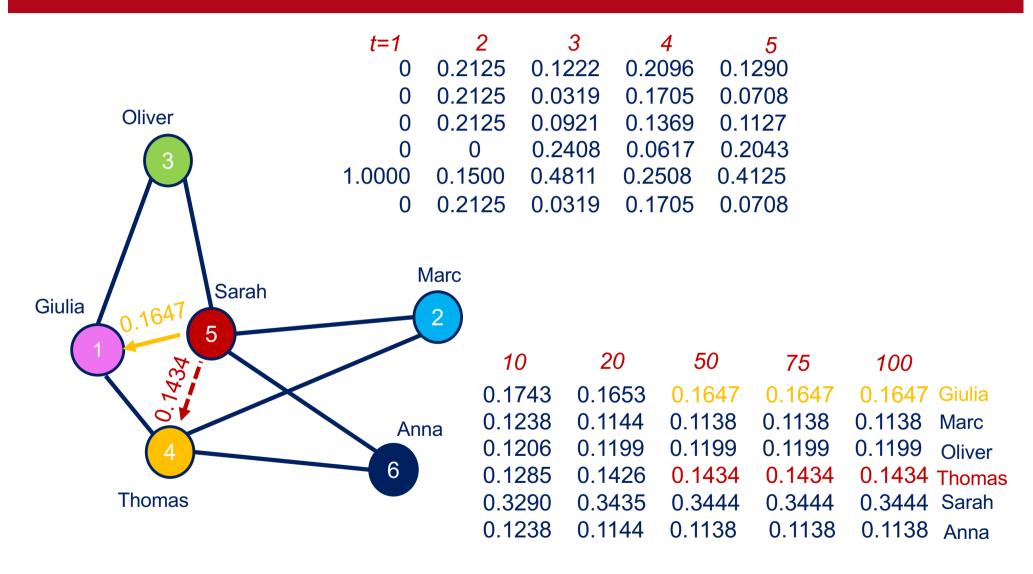
#### Result

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

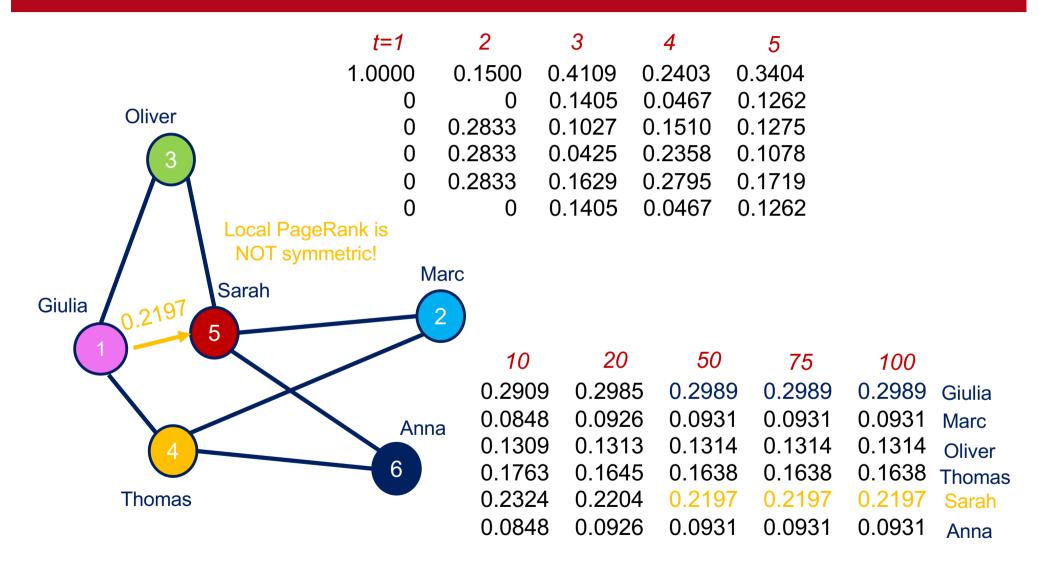




who's Sara's best friend?



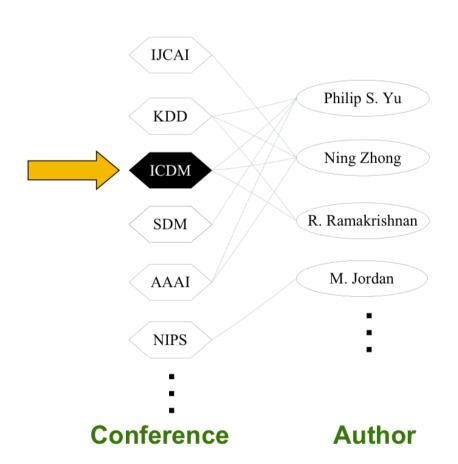
who's Giulia's best friend?

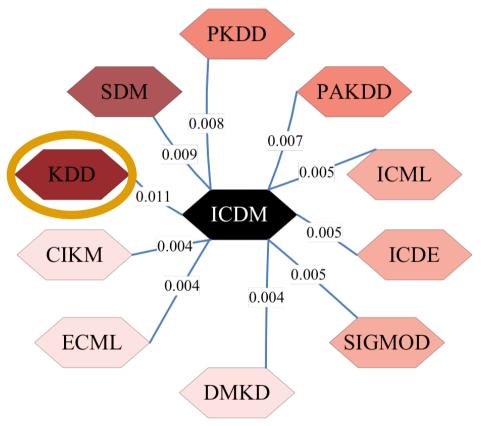




what is the most related conference to ICDM?

#### Top 10 ranking results





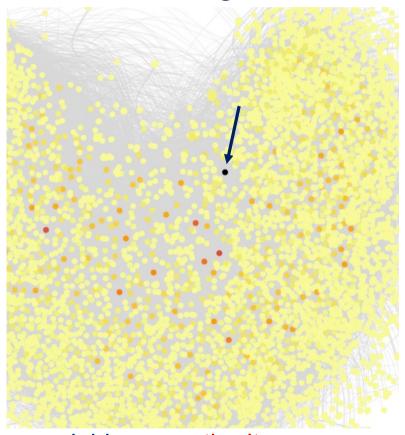
ICDM = international conf. on data mining KDD = knowledge discovery and data mining



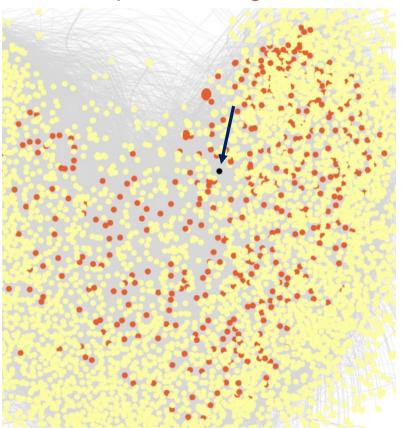
#### Local PageRank versus degree authorities

#### Local PageRank

### 1-hop out-neighbours



neighbours authority score = local node → neighbours



#### 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

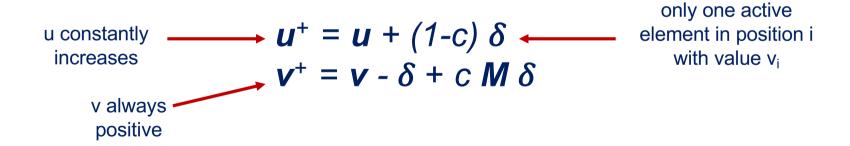
use institutional Sign In with your unipd credentials



#### Approximate PageRank algorithm

the push operation

- $\Box$  Start from u = 0 and v = q
- degree of node i
  precision sum of the degrees
- To all the nodes  $\underline{i}$  satisfying  $\underline{v_i} > \varepsilon \underline{d_i}/\underline{D}$  apply the push operation



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



### Linearity of PageRank

to build a lemma for the proof

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

□ PageRank equation  $r_q = c M r_q + (1-c) q$ 

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

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

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



$$r_{au+bv} = a r_u + b r_v$$

#### Modifying the PageRank equation

the lemma for the proof

one-step random walk

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



$$r_q = (I - c M)^{-1} (1-c) q$$

$$\Box$$
  $M r_q = (1-c) \sum (c M)^k M q$ 

$$\square$$
  $M r_q = r_{Mq}$ 



# Main property of push: $r_q = u + r_V$

- lacktriangle At starting point  $oldsymbol{u} = oldsymbol{0}$  and  $oldsymbol{v} = oldsymbol{q}$  imply  $oldsymbol{r}_q = oldsymbol{0} + oldsymbol{r}_q$
- The following steps are proved by induction

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

by linearity
$$u^{+} + r_{V+} = u + (1-c) \delta + r_{V} - r_{\delta} + c r_{M\delta}$$

$$r_{\delta} - (1-c) \delta$$

$$u^{+} + r_{V+} = u + r_{V} = r_{q}$$



# Precision guarantee: $|\mathbf{r}_q - \mathbf{u}|_1 < \varepsilon$ and the result is proved

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

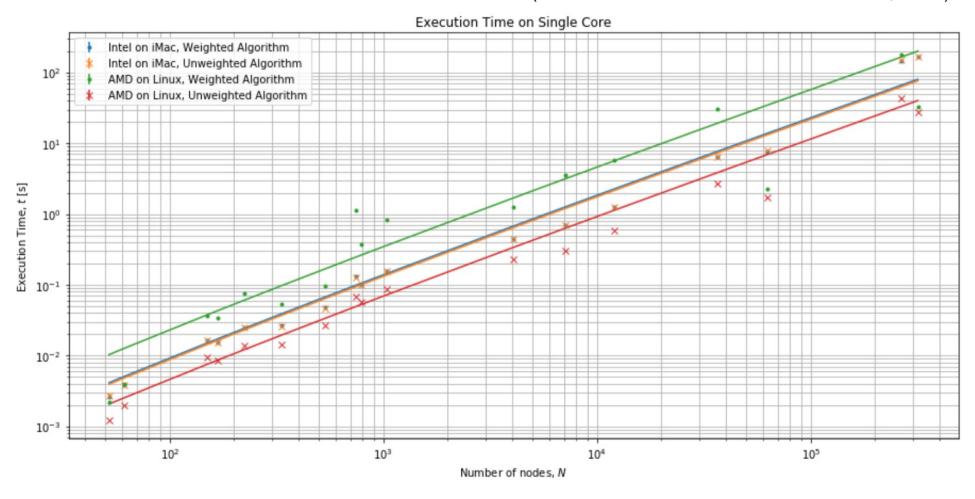
As a result 
$$|\mathbf{r}_q - \mathbf{u}|_1 = \mathbf{1}^T \mathbf{v} < \Sigma \varepsilon d_i/D = \varepsilon$$



### Scalability properties

of Local PageRank using Approximate PageRank

(Francesco Barbato & Tommaso Boccato, 2020)



#### Beware of the Lazy PageRank

which is suggested in the paper

Lazy PageRank 
$$r = a M_2 r + (1-a) q$$

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

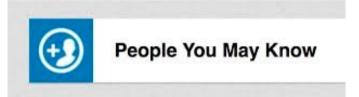
- Lazy because a fraction b of the times the surfer stays where she/he is

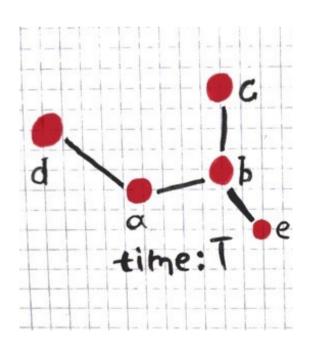
slower algorithm, as its convergence speed depends on a>c, better use c directly!



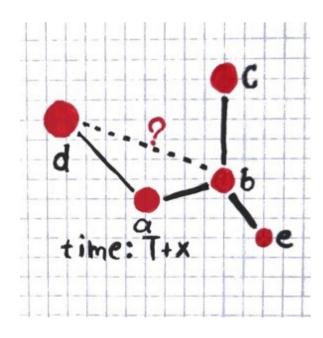
# Application #1 the link prediction task

#### Recommendation in social networks



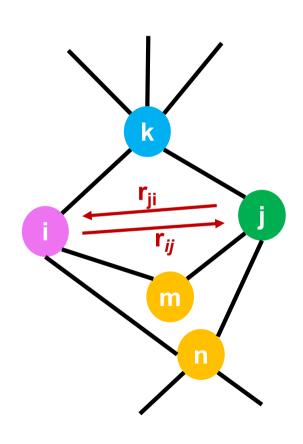


Given a graph at time
T, can we output a
ranked list of links
that are predicted to
appear in the graph
at time T+x?



#### Application #1

random walk with restart (RWR) method



Local PageRank teleportation to node 
$$i$$

$$\mathbf{r}_{i} = \mathbf{c} \ \mathbf{M} \ \mathbf{r}_{i} + (1-\mathbf{c}) \ \mathbf{\delta}_{i}$$

Likelihood of activating the link (i,j)

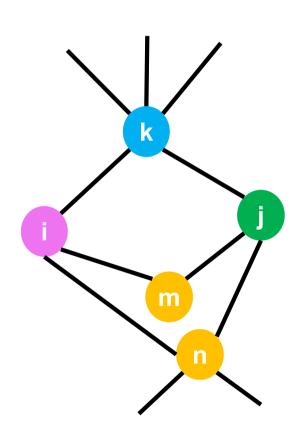
$$L_{RWR}(i,j) = r_{ij} + r_{ji}$$

Select the highest values of L<sub>RWR</sub> for recommendation pourposes



#### Application #1

the resorse allocation (RA) counterpart



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

related to a two-hop RWR

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

to have

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

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

## Application #1

performance metrics

fraction of links correctly guessed (out of 100 recomendations)

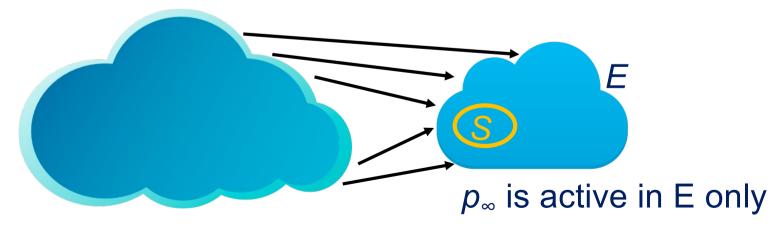
Precision	CN	RA	LP	ACT
USAir	0.59	0.64	0.61	0.49
NetScience	0.26	0.54	0.30	0.19
Power	0.11	0.08	0.13	0.08
Yeast	0.67	0.49	0.68	0.57
C.elegans	0.12	0.13	0.14	0.07
	RWR	HSM	LRW	SRW
	0.65	0.28	0.64(3)	<b>0.67</b> (3)
	0.55	0.25	0.54(2)	0.54(2)
	0.09	0.00	0.08(2)	0.11(3)
	0.52	0.84	<b>0.86</b> (3)	0.73(9)
	0.13	0.08	<b>0.14</b> (3)	<b>0.14</b> (3)

Among the best performance in social networks

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

# Application #2 TopicSpecific PageRank

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





**Tweets** 

### Application #2

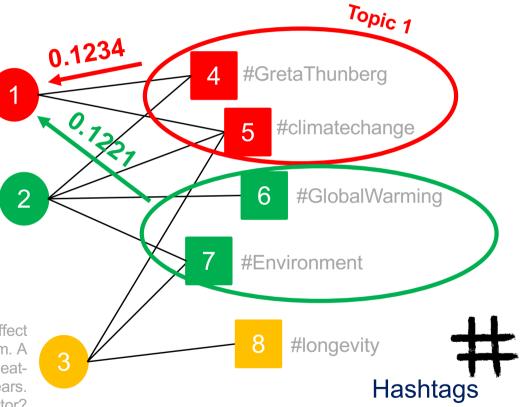
assigning documents to topics in semantic networks

#### 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 heatinduced #heartattack risk in recent years. Could #ClimateChange be a risk factor? #longevity



# Signed PageRank

modifications for signed networks



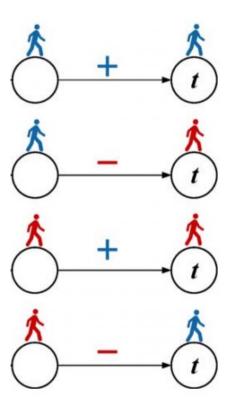
#### PageRank in signed networks

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

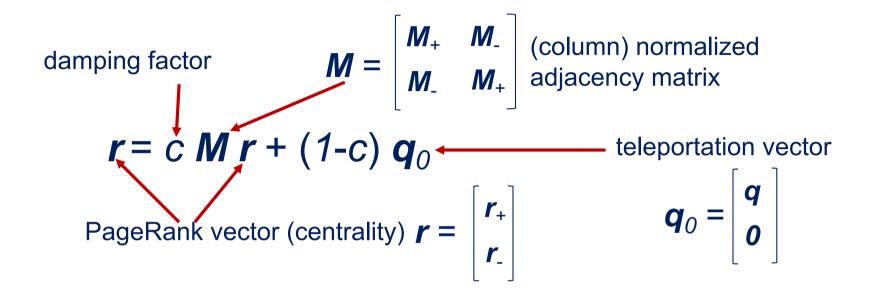
- □ Identify + (favourable) and (adversarial) paths, i.e., ranking values r<sub>+</sub> and r<sub>-</sub> for positive and negative surfers
- Extract positive  $A_+$  and negative  $A_-$  contributions to  $A = A_+ A_-$
- Normalize the absolute value, to get  $M_+$  and  $M_-$  (with normalized  $M_++M_-$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

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



#### Signed PageRank

power iteration

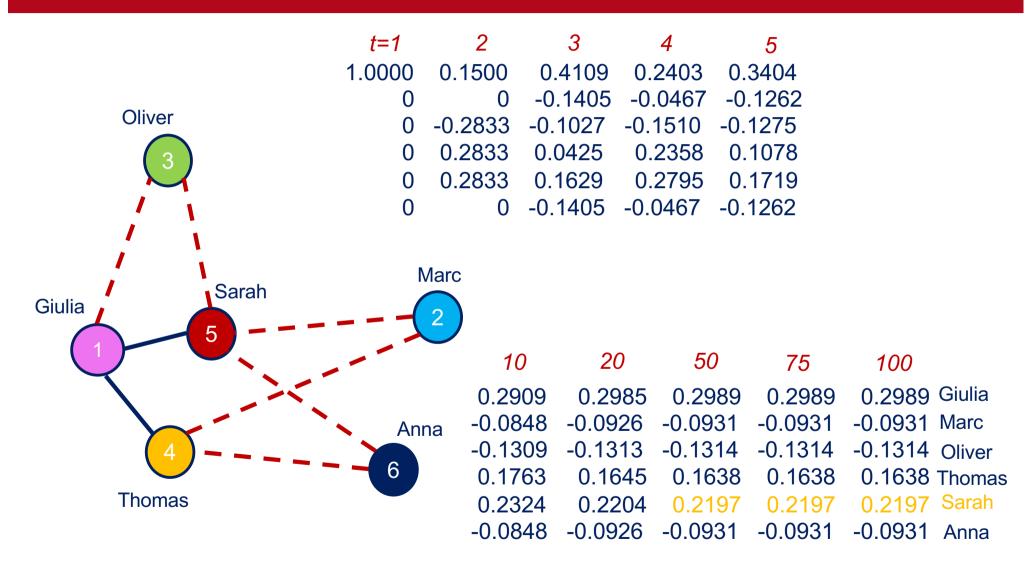


signed centrality outcome  $r_{+} = r_{+} - r_{-}$ 

$$r_{+-} = c M_{+-} r_{+-} + (1-c) q$$
 can be signed
$$M_{+-} = A \operatorname{diag}^{-1}(|A|^{T}1)$$



who's Giulia's best friend?

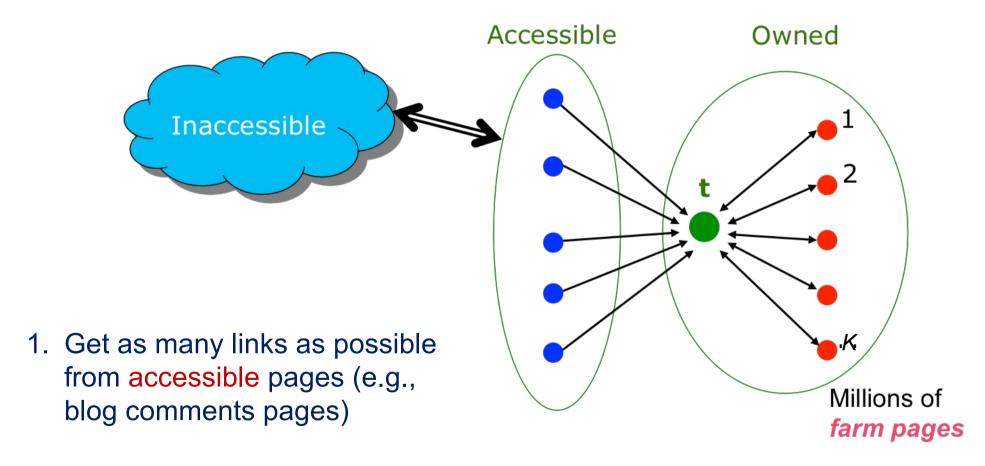


# Preventing spamming

on the role of the teleport vector

#### Spam farm

how to boost PageRank for a web page



2. Construct link farm to get a PageRank multiplier effect



# Google bombs in 2004 US elections



Web Images Groups News Froogle Local more »

miserable failure

Search

Advanced Search Preferences

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 - 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/ - 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 - 31k - Cached - Similar pages

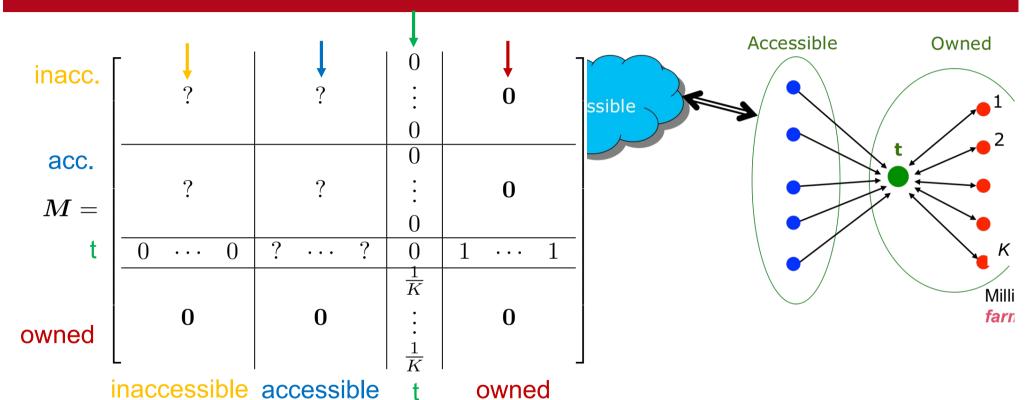
#### Google's (and Inktomi's) Miserable Failure

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 - 45k - Sep 1, 2005 - Cached - Similar pages



# PageRank analysis of spam farms



ranking due to accessible pages

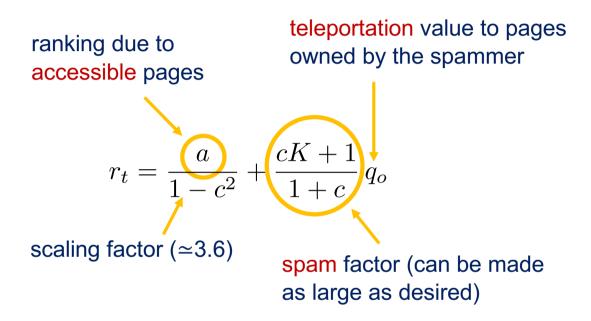
teleportation value to pages owned by the spammer

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

$$r_{t} = a + cKr_{o} + (1-c)q_{o}$$
 
$$r_{o} = c\frac{1}{K}r_{t} + (1-c)q_{o}$$
 113



# PageRank outcome of spam farms



### 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



## Row-normalized PageRank

an overview

M 1 = 1

PageRank equation 
$$r = c M r + (1-c) q$$
 row-normalized  $M = \text{diag}^{-1}(\mathbf{d}) A, d = A 1$  row-normalized  $M = 1$ 

Markov chain 
$$p_{t+1} = c M p_t + (1-c) q$$

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

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

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

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

same properties of column-normalized PageRank:

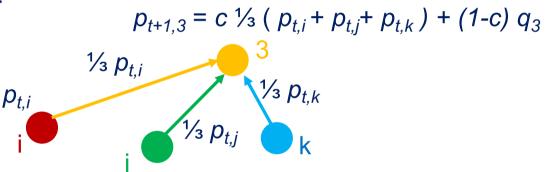
- M₁ has one eigenvalue equal to 1
- The remaining eigenvalues satisfy  $|\lambda| \leq c$



### Row-normalized PageRank

interpreting its action

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



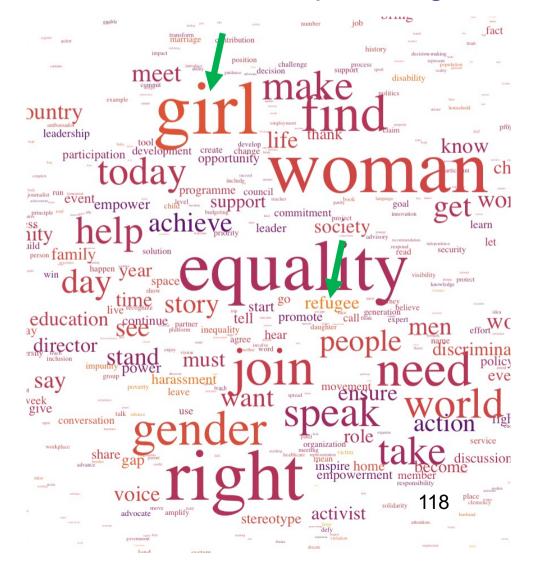
### Semantic network example

agency = action and goal orientation, sense of which is necessary for people to attempt social change

#### q values of agency (in colour)

# ountry create change opportunity participation development create know child Support director say activist

#### r values after spreading



# Takeaways for PageRank centrality

- 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 centrality

hubs and authorities



#### 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

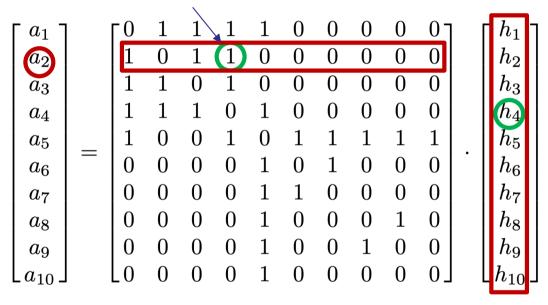
We deprecate its use



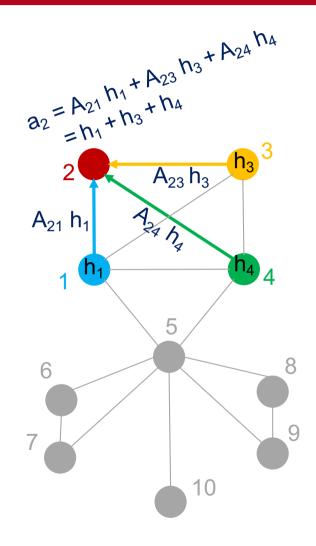
#### HITS equations

authorities score









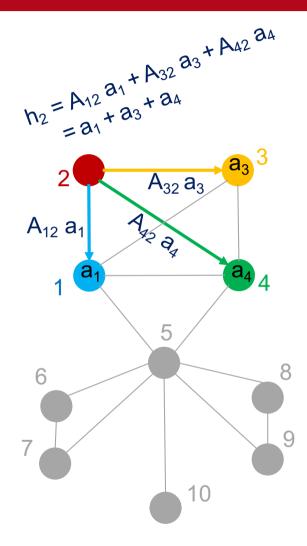


#### HITS equations

hubs score



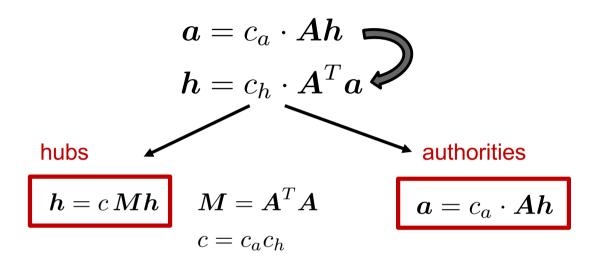
$$h = A^T a$$

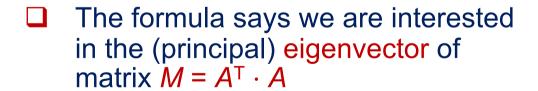




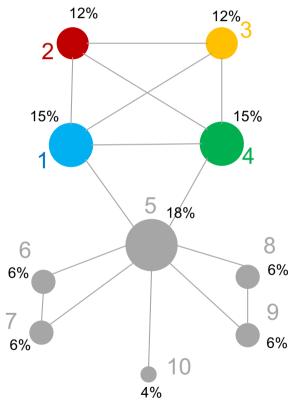
#### HITS equations

hubs and authorities











# Power iteration method for HITS

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

1. Let the time go by
$$a_{t+1} = M a_t$$
product by a sparse matrix (twice)  $M = A A^T$ 

2. Keep normalizing (divide  $a_{t+1}$  by the sum of elements)

3. Stop when *a* converges (few iterations?)

# Convergence properties

- $\square$   $\lambda_1$  largest eigenvalue of M
- $\square$   $\lambda_2$  second largest eigenvalue of **M**
- □ Triang. inequality  $\|\boldsymbol{a}_{t}-\boldsymbol{a}_{t+1}\|_{2} \leq 2\sqrt{N} \cdot (\lambda_{2}/\lambda_{1})^{t}$

#### Worst case result:

- □ Precision  $\varepsilon$  implies:  $\|\mathbf{a}_{t} \mathbf{a}_{t+1}\|_{2} < \varepsilon$
- □ Iterations required:  $t = [\ln(2/\epsilon) + \frac{1}{2}\ln(N)] / \ln(\lambda_1/\lambda_2)$

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

$$\frac{\ln(N)}{\ln(\lambda_1/\lambda_2)}$$

 $N = 10^9 \rightarrow 10.3$ 

# Eigenvector and Kats centralities

other (less interesting) alternatives to PageRank

 $r = (I - c A)^{-1} 1$ 

 $= \Sigma (c A)^k 1$ 

#### Eigenvector and Kats centralities

an overview

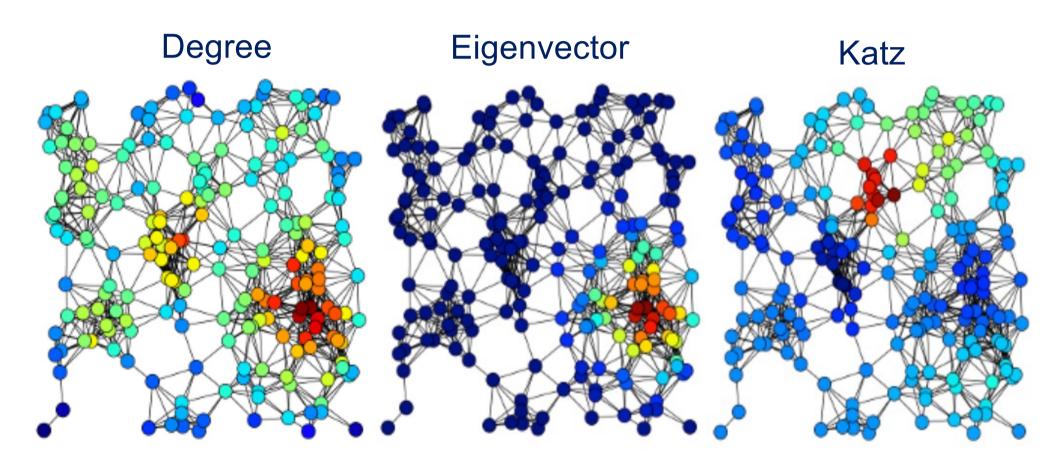
	with constant term	without constant term
red	PageRank	Degree
normalized	r = c M r + (1-c) q	r = M r
ized	Katz	Eigenvector
unnormalized	r = c A r + 1	r = c <b>A</b> r
		The charge of

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



## Eigenvector and Kats centralities

their graphical interpretation





# Closeness and Harmonic centralities

importance of nodes as spreaders of information

## Closeness centrality

a definition

### 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, [1][2] that is:

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

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

ser it is to

If the Rationale: the node which is the one which which the one which are ach, the one which is the pest for spreading easiest to reach, spreading information information

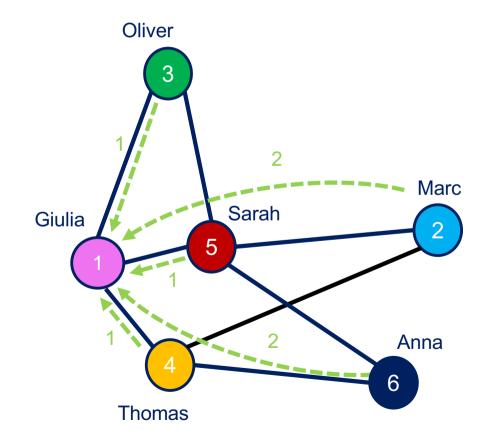


#### An example

on how to calculate closeness centrality

count the lengths of the shortest paths leading to Giulia





#### 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

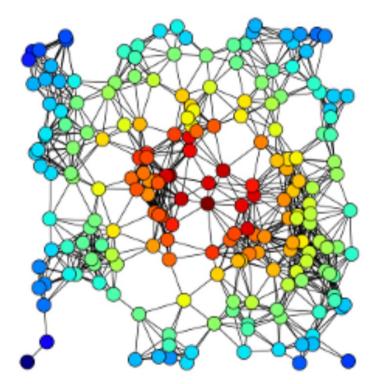
$$C(Giulia) = 1/7$$
  
= 0.1429



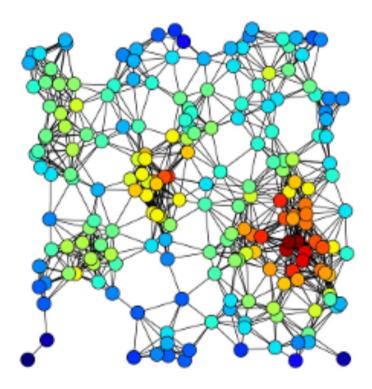
## Closeness versus degree centrality

a graphical interpretation

#### Closeness



#### Degree







## Harmonic centrality

a definition

#### 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  $1/\infty=0$ :

$$H(x) = \sum_{y 
eq x} rac{1}{d(y,x)}.$$

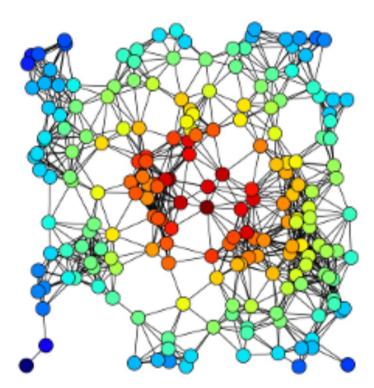
The most natural modification of Bavelas's definition of closeness is following the general principle proposed by Marchiori and Latora (2000)<sup>[3]</sup> that in graphs with infinite distances the harmonic mean behaves better than the arithmetic mean. Indeed, Bavelas's closeness can be described as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



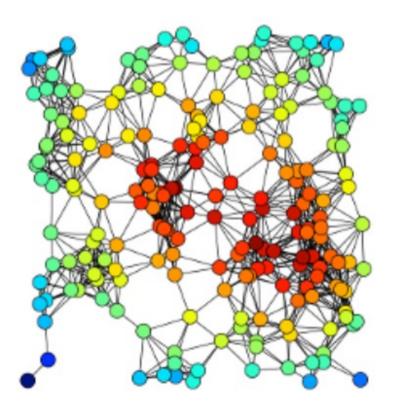
## Closeness versus harmonic centrality

a graphical interpretation

#### Closeness



#### Harmonic





# Betweenness centrality

importance of nodes as bridges or brokers

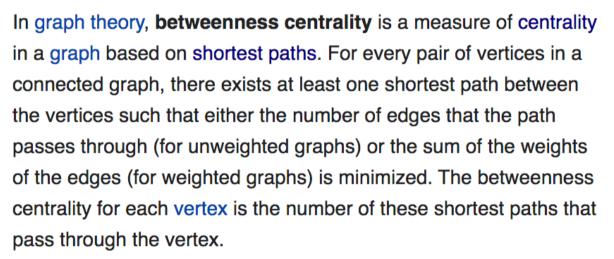


### Betweenness centrality

a definition

#### Betweenness centrality

From Wikipedia, the free encyclopedia



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.



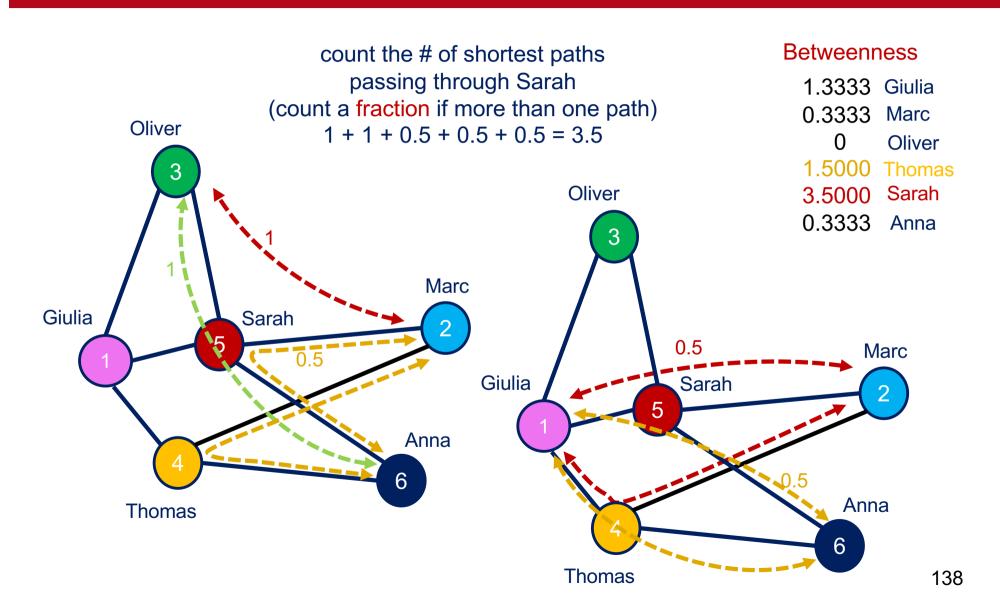
Rationale: the node which takes

you elsewhere broker)

you elsewhere bridge, broker)

#### An example

on how to calculate betweenness centrality



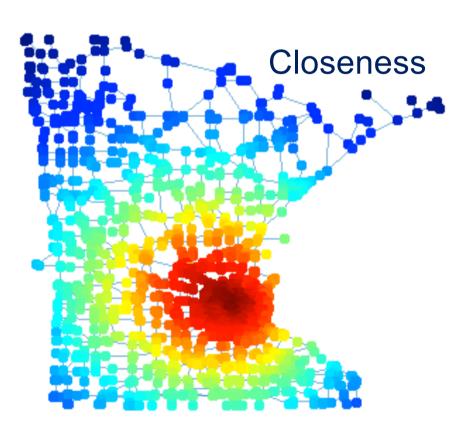


#### Closeness vs betweenness centrality

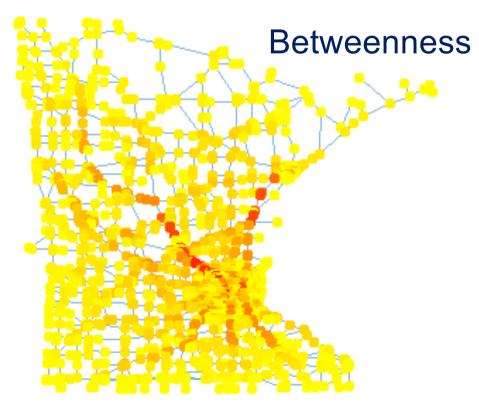
a graphical interpretation

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)



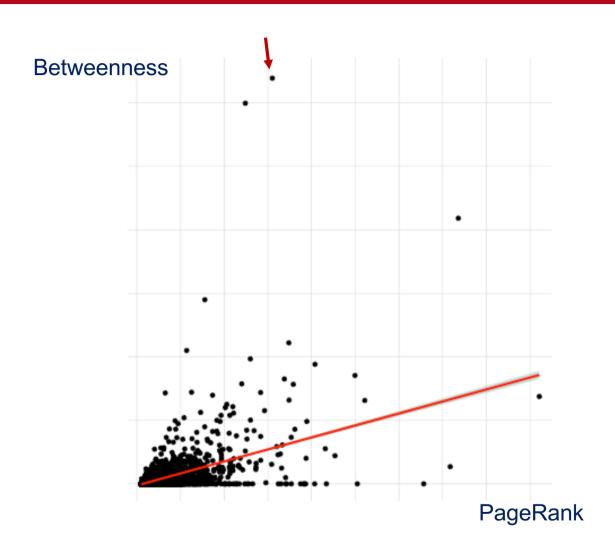
## Betweenness vs PageRank centrality

wiki vote network



## Betweenness vs PageRank centrality

a correlation view



# Clustering coefficient

how tightly linked is the network locally



## Clustering coefficient

a definition

#### Local clustering coefficient [edit]

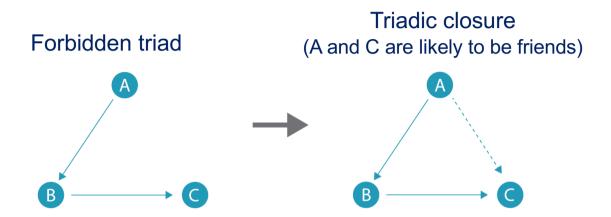


The **local clustering coefficient** of a vertex (node) in a graph quantifies how close its neighbours are to being a clique (complete graph). Duncan J. Watts and Steven Strogatz introduced the measure in 1998 to determine whether a graph is a small-world network.

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

#### Triadic closure

in social networks

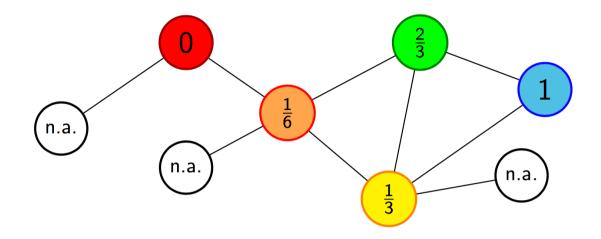


#### Triadic closure

- □ 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

$$C_i = \frac{1}{|\mathcal{N}_i|(|\mathcal{N}_i| - 1)} \sum_{\substack{(j,k) \in \mathcal{N}_i^2 \\ i \neq k}} \operatorname{tc}_{i,j,k}$$
 equal to diag(A³)

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

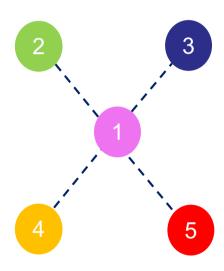


## Local clustering coefficient

examples

not connected
neighbourhood

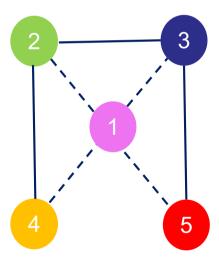
$$< C > = 0$$



$$C_1 = 0$$

weakly connected neighbourhood

$$< C > = 0.766$$



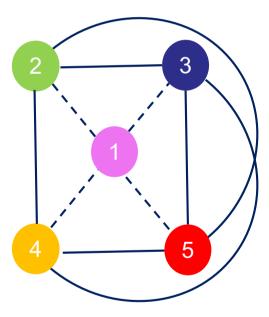
$$C_1 = \frac{1}{2} = \frac{3}{4x3/2}$$

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

$$C_4 = C_5 = 1$$

strongly connected neighbourhood

$$< C > = 1$$



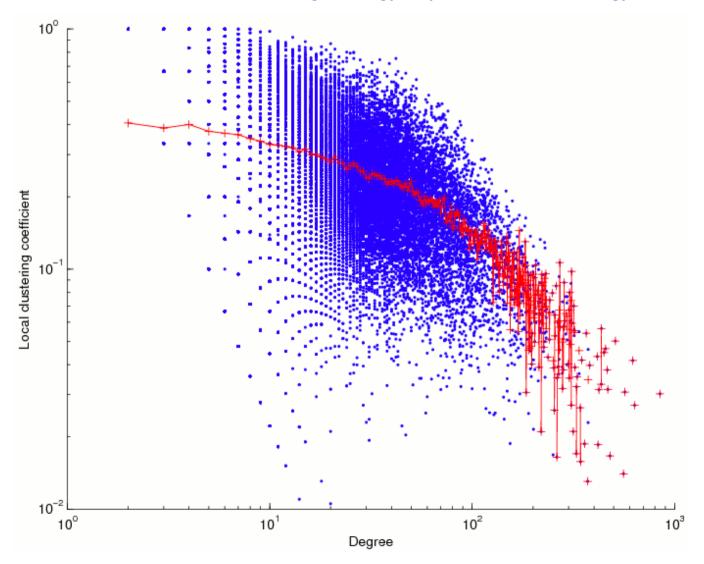
$$C_1 = 1 = 6 / (4x3/2)$$



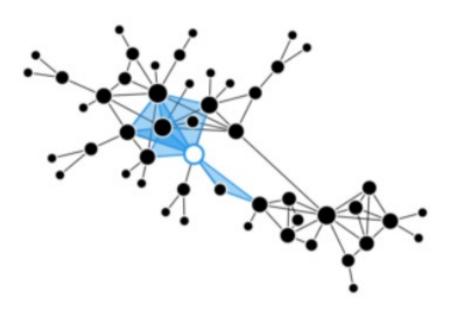
### Clustering coeff. vs degree centrality

a correlation view

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



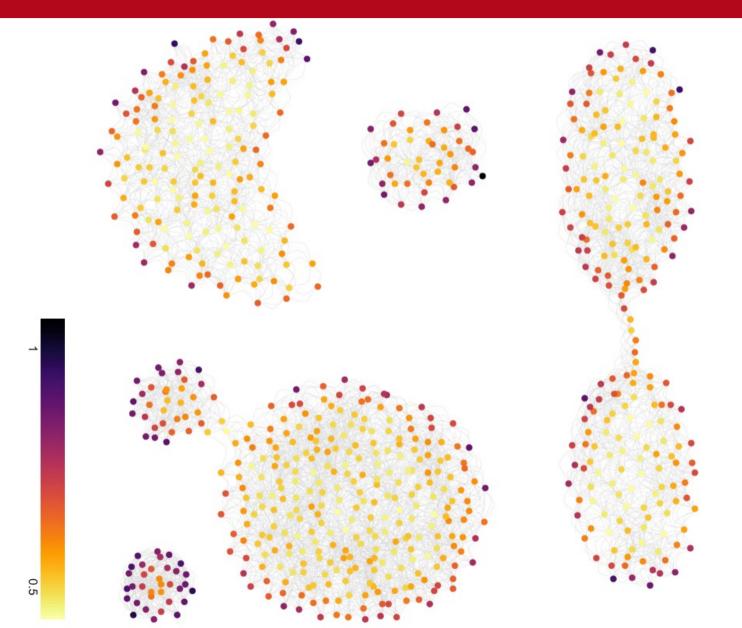
when person has many friends, these friends have less edges among them, which is to be expected since a person with many friends is likely to have friends from more diverse communities, and a paper getting cited many times is likely to be cited by papers from more diverse areas



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

## Visual example

149





for Closeness, Betwenness and Clustering coefficient

- 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

# Takeaways on centrality measures

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



#### More on the meaning

https://reticular.hypotheses.org/1745

