

Network Science

#5 PageRank centrality

© 2020 T. Erseghe



What is centrality?

Centrality

From Wikipedia, the free encyclopedia



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.^[1]

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

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

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

Main article: [PageRank](#)

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

Main article: [Betweenness centrality](#)

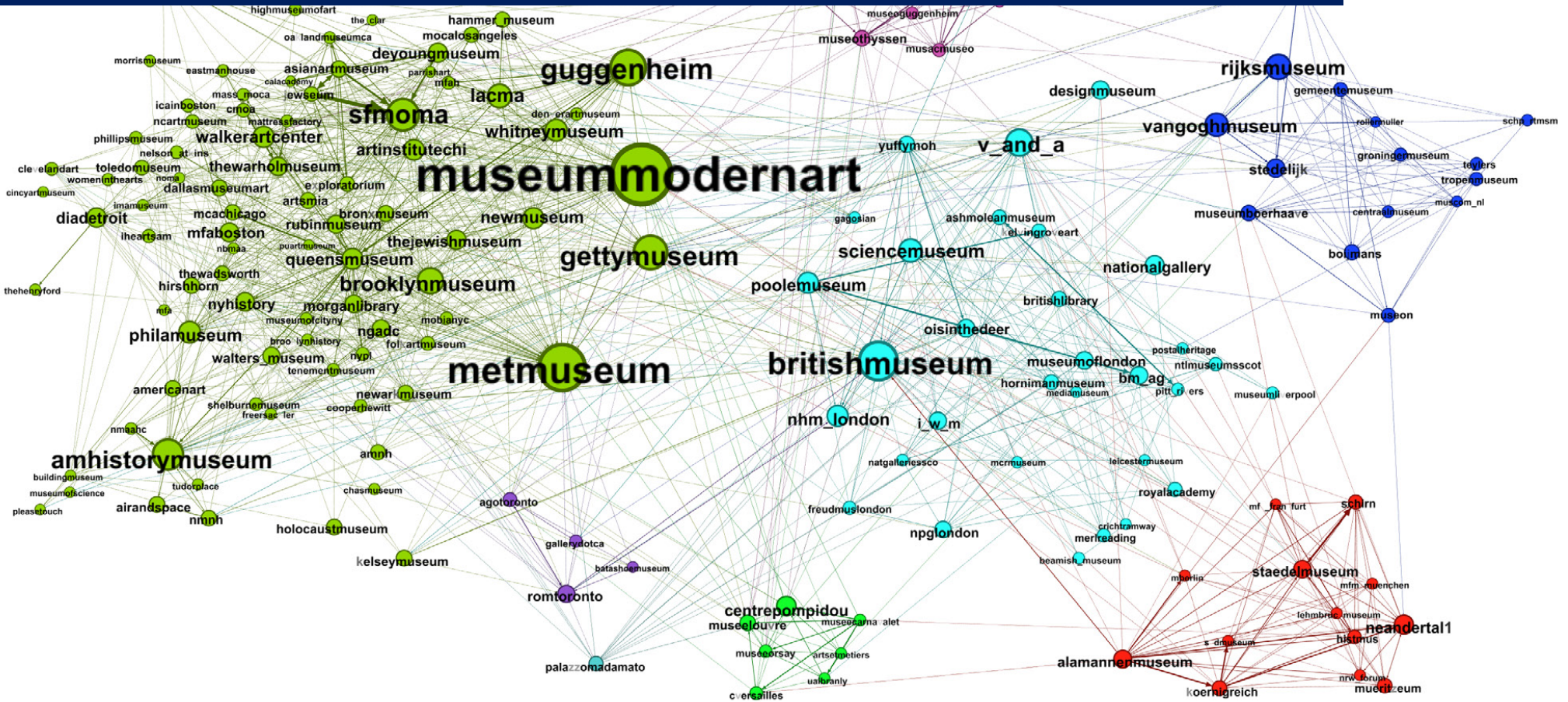
[Eigenvector centrality](#) [\[edit \]](#)

Main article: [Eigenvector centrality](#)

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

Main article: [Closeness centrality](#)

How to rank nodes in a network?



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

The solution comes from the web 
VIME.

How to organize the web?

Idea: 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



Two approaches



PageRank

Page, Brin, Motwani, Winograd
1999

«The PageRank citation ranking:
bringing order to the web»
Stanford InfoLab

HITS – hubs and authorities

Kleinberg, J.M.
1999

«Authoritative sources in a
hyperlinked environment»
Journal of the ACM

Conceptually similar

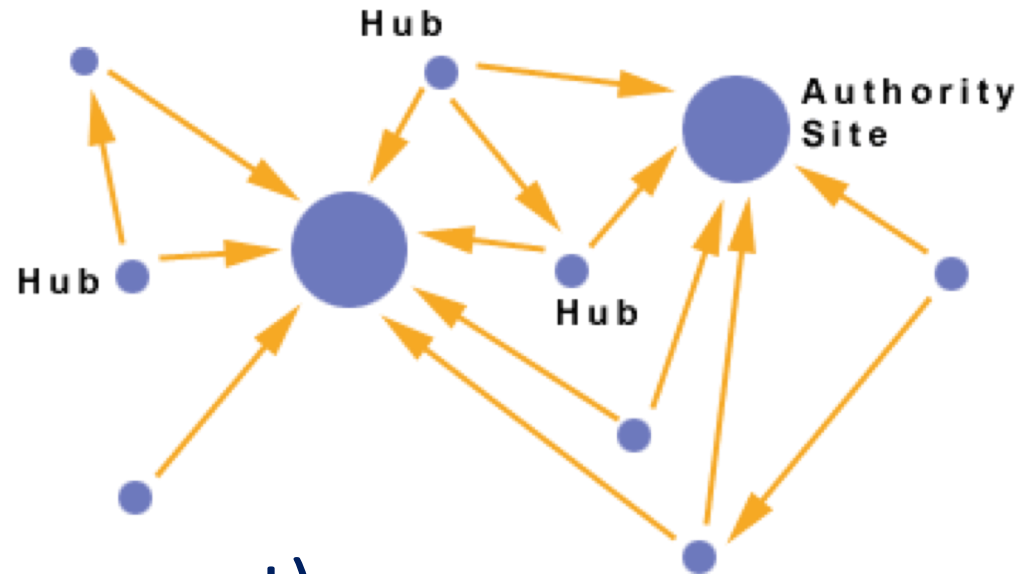
HITS centrality

Hiperlink induced topic search (HITS)

Two classes of nodes:

❑ **Authorities** (quality as a content provider)

nodes that contain useful information, or having a high number of edges pointing to them (e.g., course homepages)



❑ **Hubs** (quality as an expert)

trustworthy nodes, or nodes that link to many authorities (e.g., course bulletin)

authority or hub?

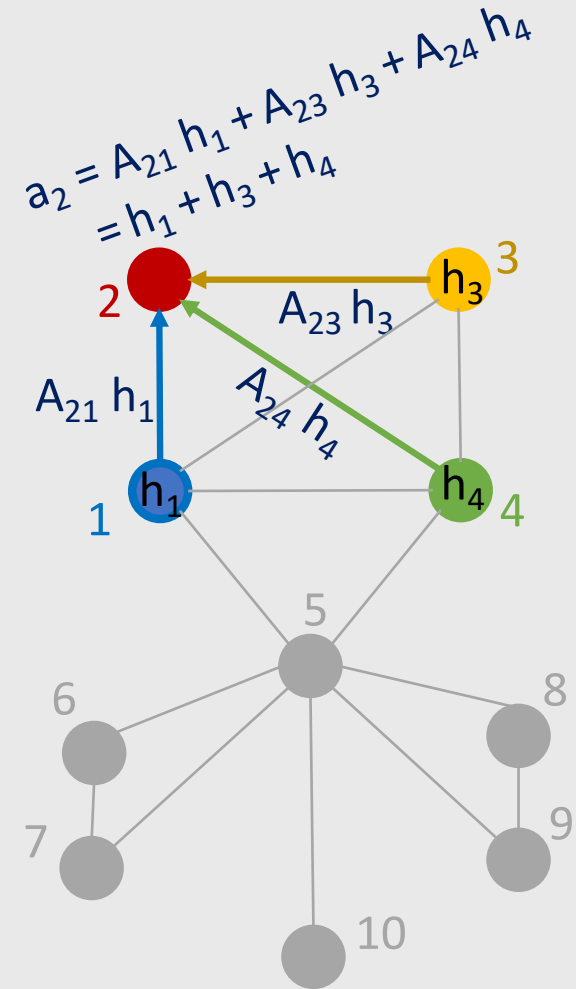


HITS equations – authority score

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

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

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



HITS equations – hub score

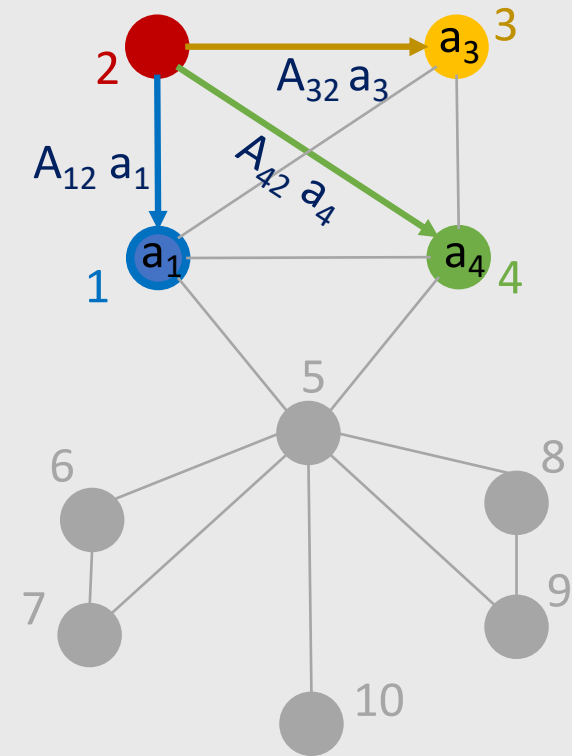
$A_{3,2}$ = weight of connection 2 \rightarrow 3

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

$$h = A^T a$$

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

$$= a_1 + a_3 + a_4$$



HITS equations

$$a = c_a \cdot Ah$$

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

hubs

$$h = c M h$$

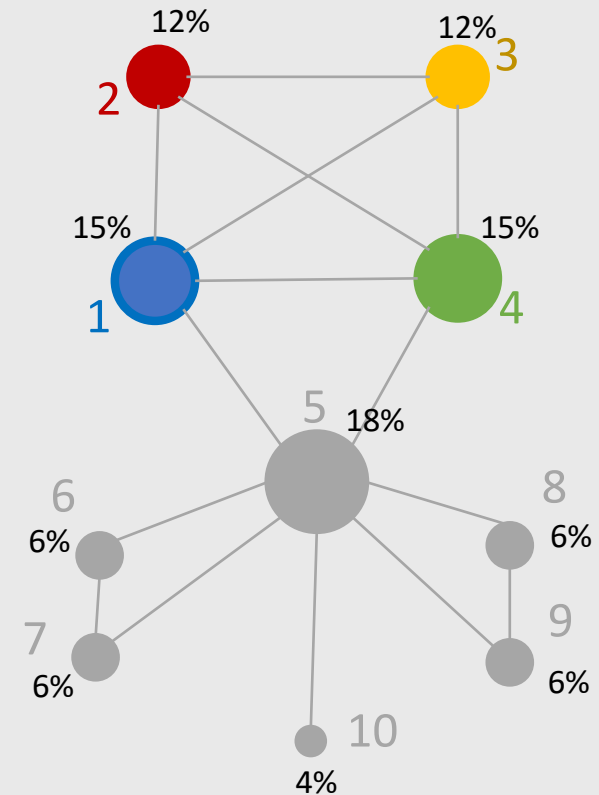
$$M = A^T A$$

$$c = c_a c_h$$

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



Power iteration method

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

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

1. Let the time go by

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

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

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

Power iteration method

Main **convergence** property:

- $\|\mathbf{a}_t - \mathbf{a}_\infty\|_2 \leq \sqrt{N} \cdot (\lambda_2/\lambda_1)^t$
- λ_1 largest **eigenvalue** of M
- λ_2 second largest eigenvalue of M
- Triang. inequality ensures $\|\mathbf{a}_t - \mathbf{a}_{t+1}\|_2 \leq 2\sqrt{N} \cdot (\lambda_2/\lambda_1)^t$

Worst case result:

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

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

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

$N = 10^9 \rightarrow 10.3$

slow if λ_2 close to λ_1

Application example – The news

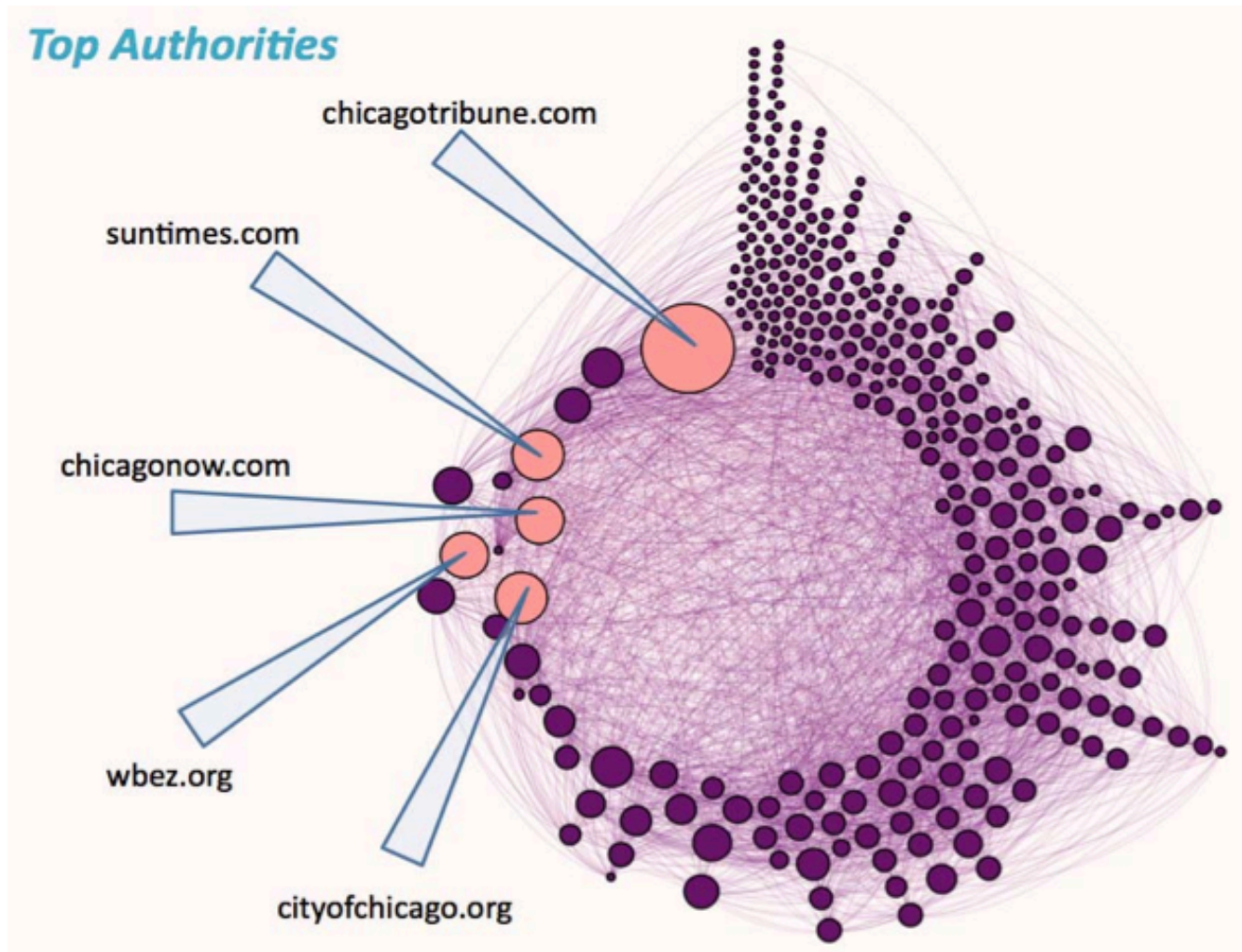
Sonderman [2012]

«Study: Smaller news websites depend more on social media for traffic than larger sites,»

- ❑ Examined links between 301 news websites, for a two-week period
- ❑ 23 percent of all their referrals were from social media
- ❑ Small websites got more than half their referrals from social media, while the large sites got only about 19 percent from social

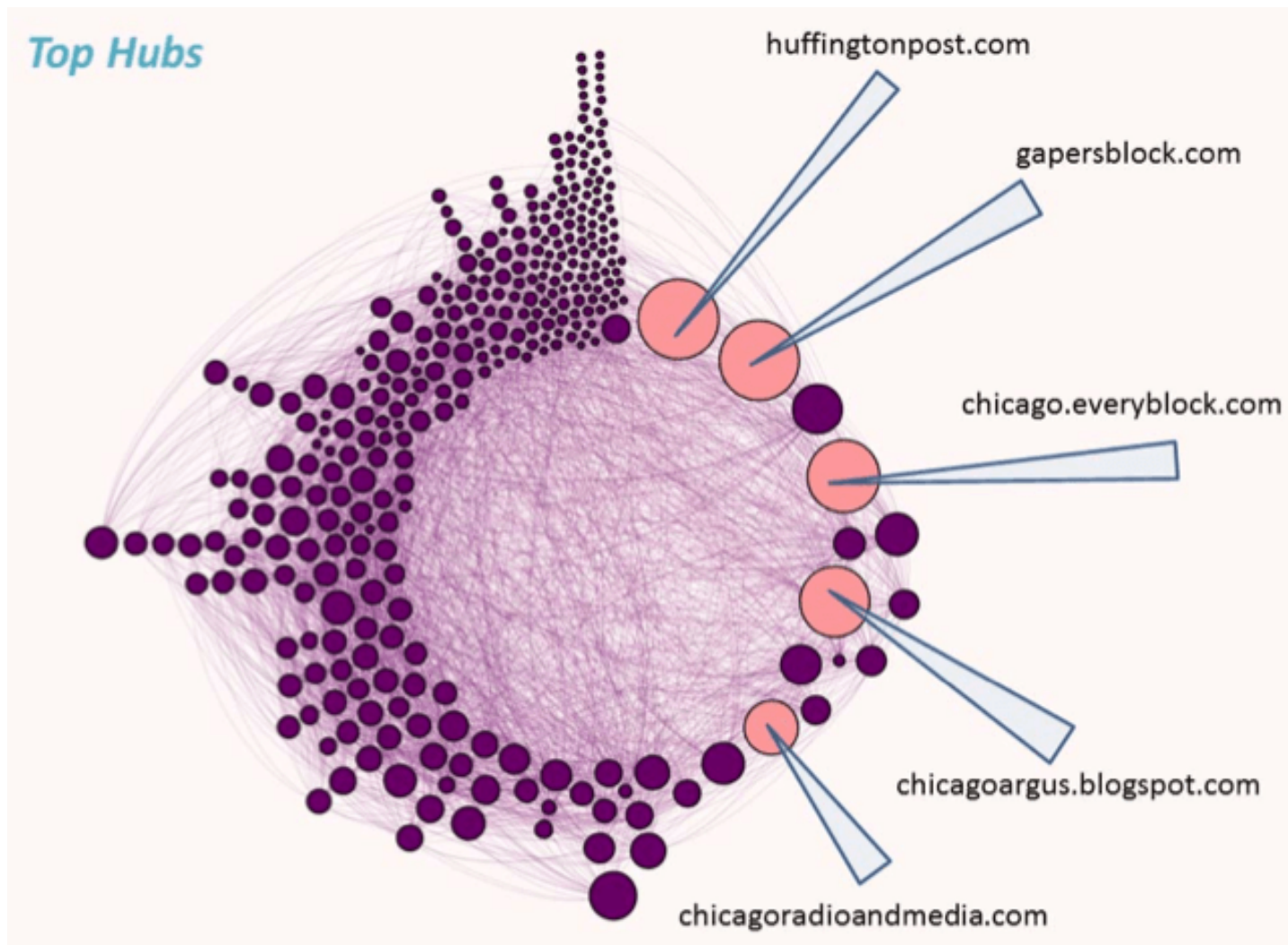


Authorities



□ The legacy media brands

Hubs



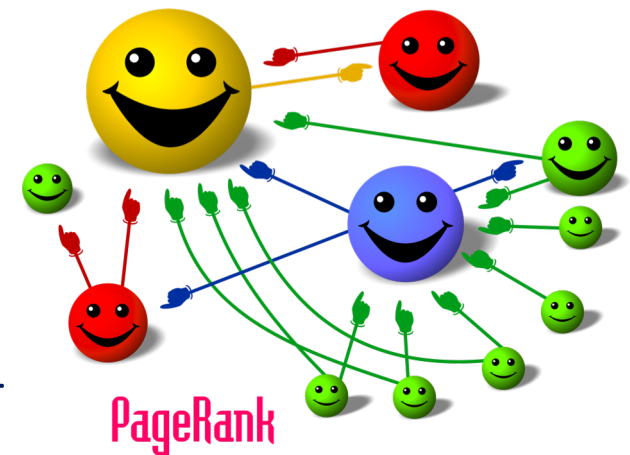
□ A different set of leading websites

PageRank centrality

PageRank

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
- ❑ Same ideas as HITS authorities
- ❑ Can be extended to hubs by using A^T

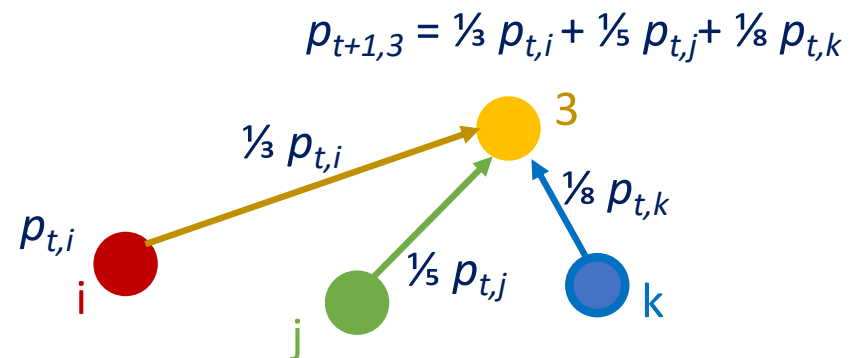
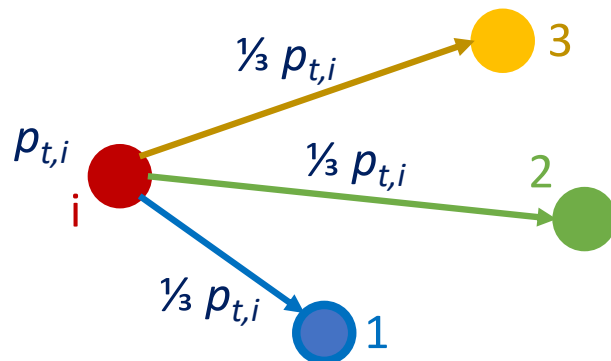


PageRank

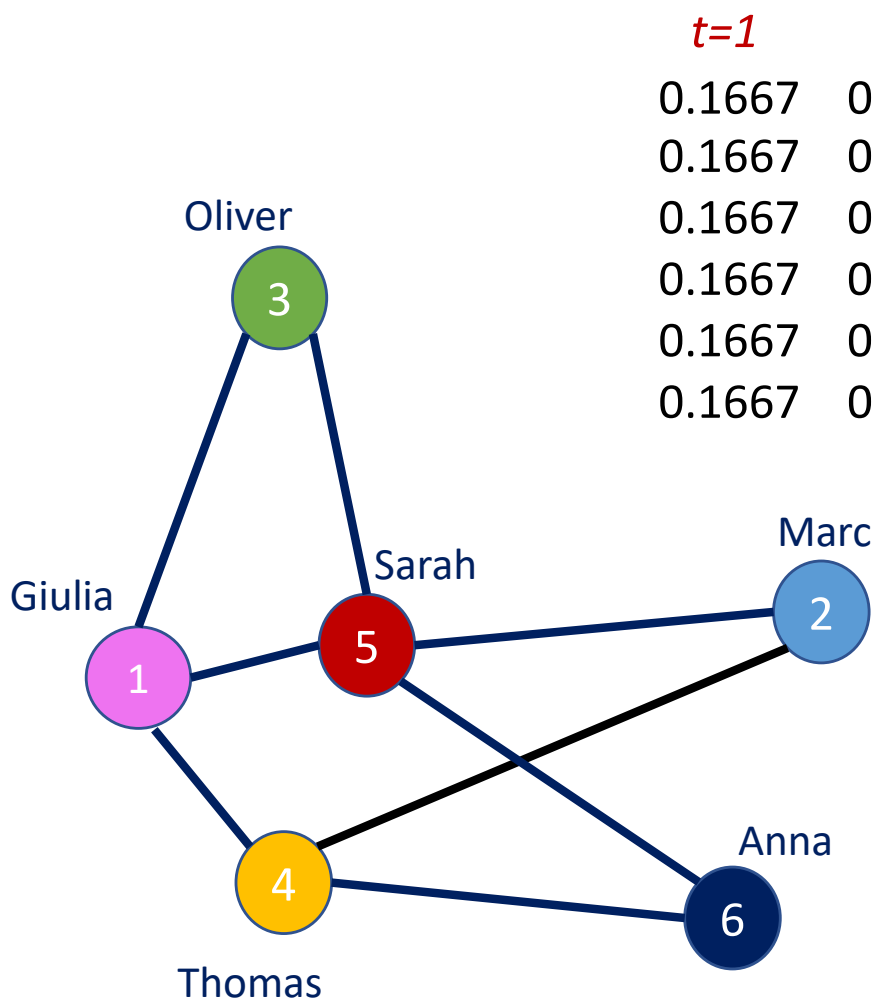


The rationale behind PageRank

- ❑ Random walk
- ❑ 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** (authority score)



Example



	<i>t=1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
Oliver	0.1667	0.1806	0.1991	0.1723	0.2025
Giulia	0.1667	0.0972	0.1505	0.1040	0.1436
Thomas	0.1667	0.0972	0.1366	0.1179	0.1287
Sarah	0.1667	0.2222	0.1574	0.2168	0.1614
Marc	0.1667	0.3056	0.2060	0.2851	0.2203
Anna	0.1667	0.0972	0.1505	0.1040	0.1436

Equal to
(normalized)
degree
centrality in
undirected
networks !!!

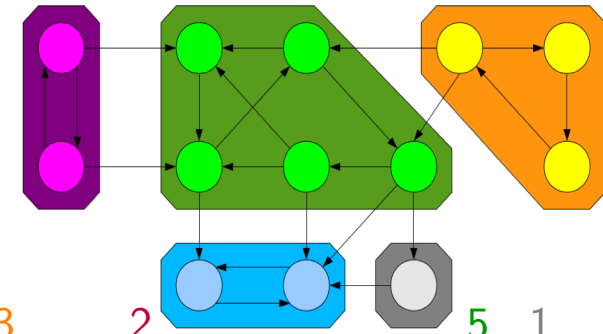


	<i>10</i>	<i>20</i>	<i>50</i>	<i>75</i>	<i>100</i>	
Giulia	0.1783	0.1848	0.1874	0.1875	0.1875	Giulia
Marc	0.1153	0.1222	0.1249	0.1250	0.1250	Marc
Oliver	0.1242	0.1248	0.1250	0.1250	0.1250	Oliver
Thomas	0.2020	0.1917	0.1876	0.1875	0.1875	Thomas
Sarah	0.2649	0.2543	0.2501	0.2500	0.2500	Sarah
Anna	0.1153	0.1222	0.1249	0.1250	0.1250	Anna

PageRank

Markov chain interpretation

- ❑ $p_{t+1} = M p_t$
- ❑ p_t stochastic vector (positive entries which sum up to 1)
- ❑ M normalized adjacency matrix (column stochastic)
- ❑ $M = A \text{diag}^{-1}(d)$
- ❑ $d = A^T \mathbf{1}$ output degree vector
- ❑ $p_\infty = M p_\infty$ converges to an eigenvector of M (with eigenvalue 1)



$M =$

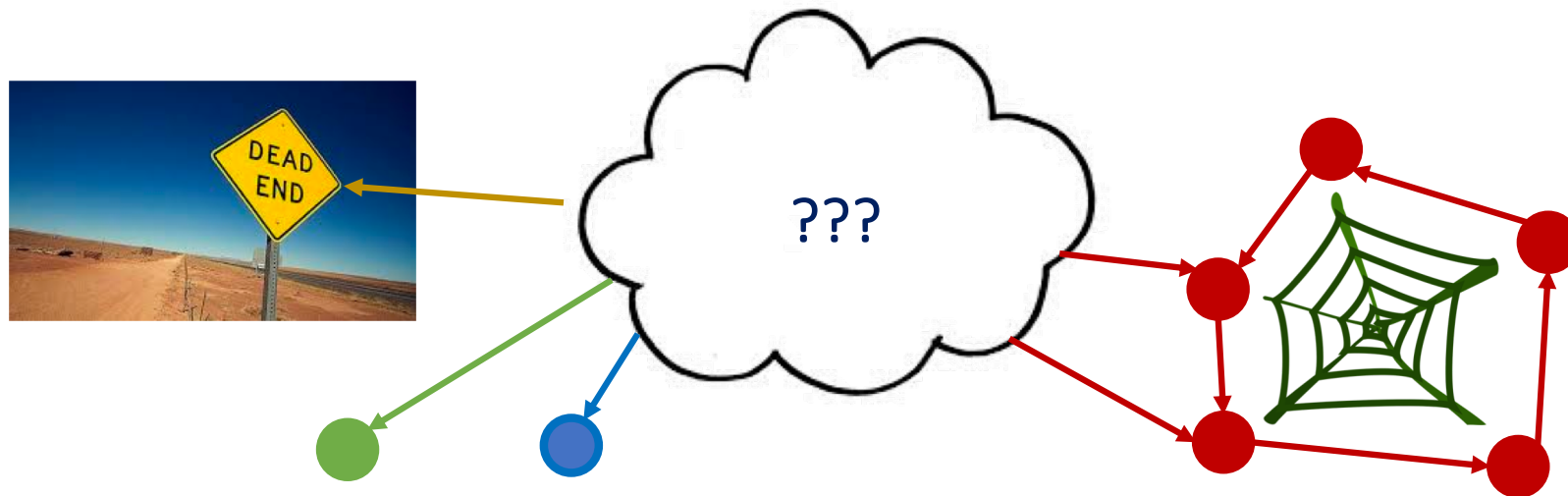
	3	2	5	1	2
1	$\frac{1}{3}$	1			
1		$\frac{1}{2}$			
1	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{2}$
1	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{1}{2}$
1		$\frac{1}{2}$	1	$\frac{1}{3}$	
1			$\frac{1}{3}$	$\frac{1}{3}$	0
1			$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{2}$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

↑ ↑
columns sum to 1

PageRank

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

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



PageRank with restart

- dead ends
- no dead ends
- normalization
- no spider traps

original adjacency matrix
(can be fractional)

A_0

teleportation vector

$$A = A_0 + b e^T$$

indicating vector
of dead ends

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

$$M_1 = c M + (1-c) q \mathbf{1}^T$$

equivalent formulation
matrix is no more sparse

$$p_{t+1} = M_1 p_t$$

PageRank with restart

PageRank
equation

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

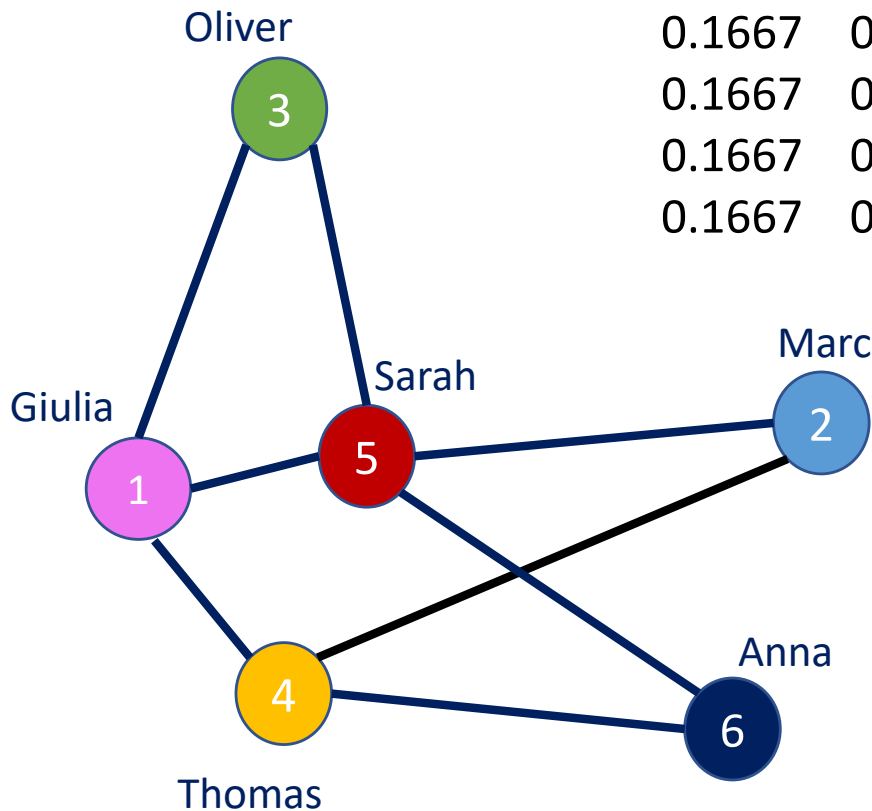
damping factor

(column) normalized adjacency matrix

teleportation vector

PageRank vector (centrality)

Example (cont'd)



	<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
Giulia	0.1667	0.2847	0.2128	0.2614	0.2276
Giulia	0.1667	0.1076	0.1461	0.1176	0.1382

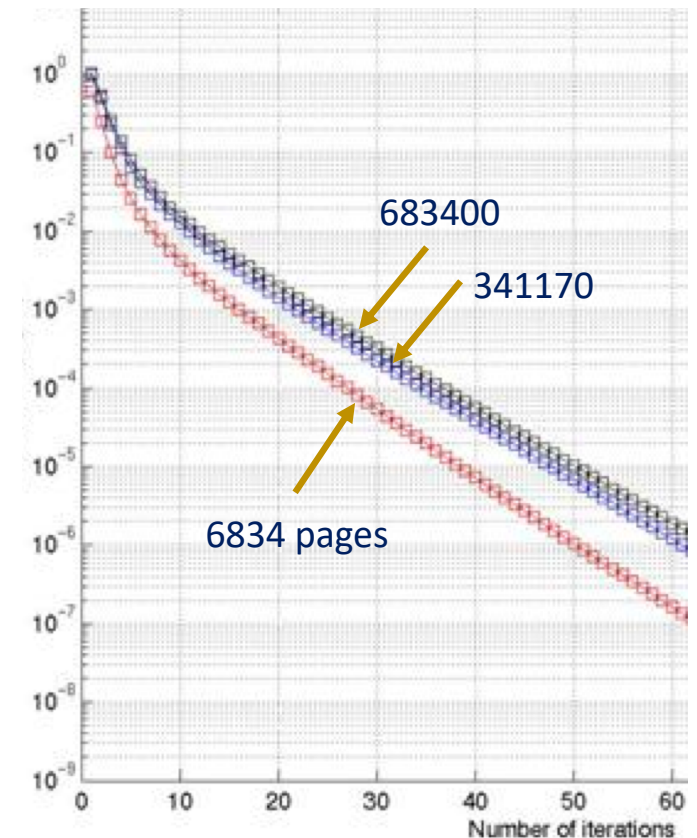
not anymore identical to degree centrality !!!



	<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

Properties

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



Power iteration method

Main **convergence** property:

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

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

Complexity considerations:

□ Precision ε at: $\|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2 < \varepsilon$

□ Iterations: $t = \lceil \ln(2/\varepsilon) + \ln(K) \rceil / \ln(1/c)$

↑
precision $10^{-3} \rightarrow 7.6$

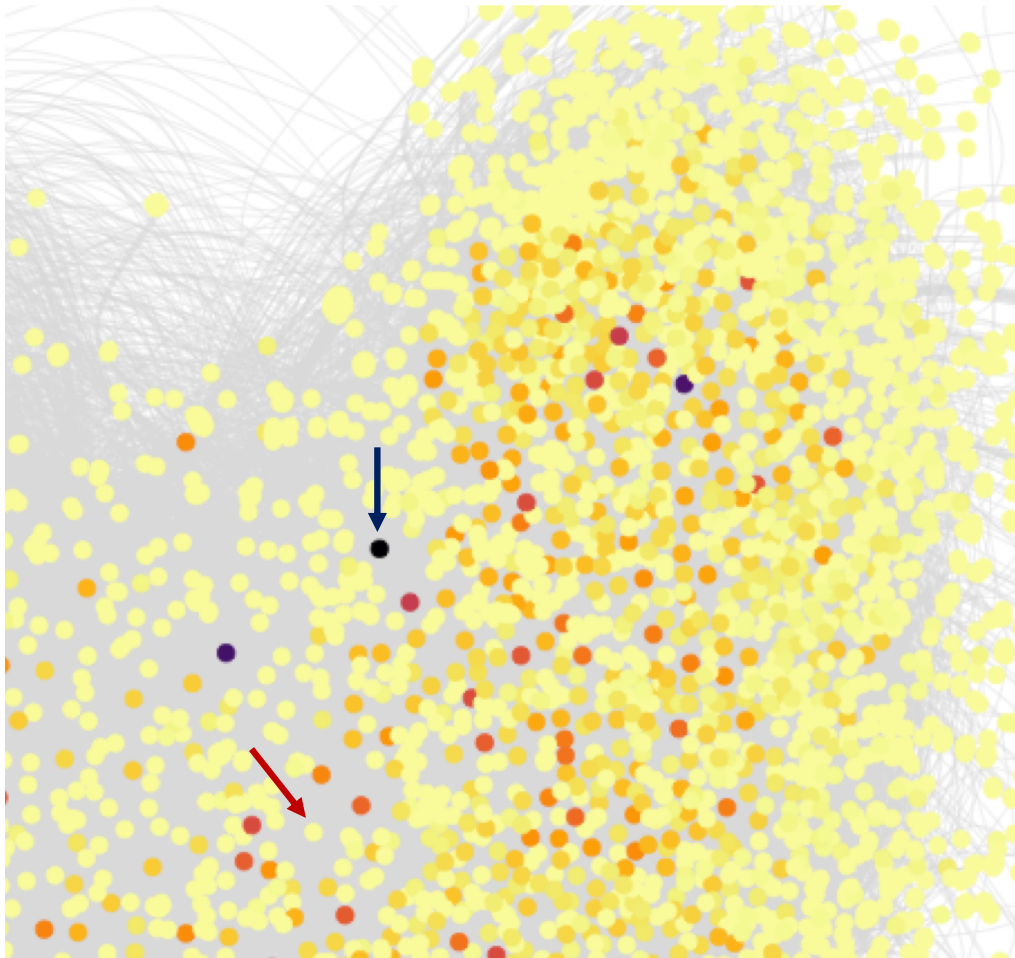
↑
can be proportional to
 $\ln(N) \rightarrow$ slow algorithm

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

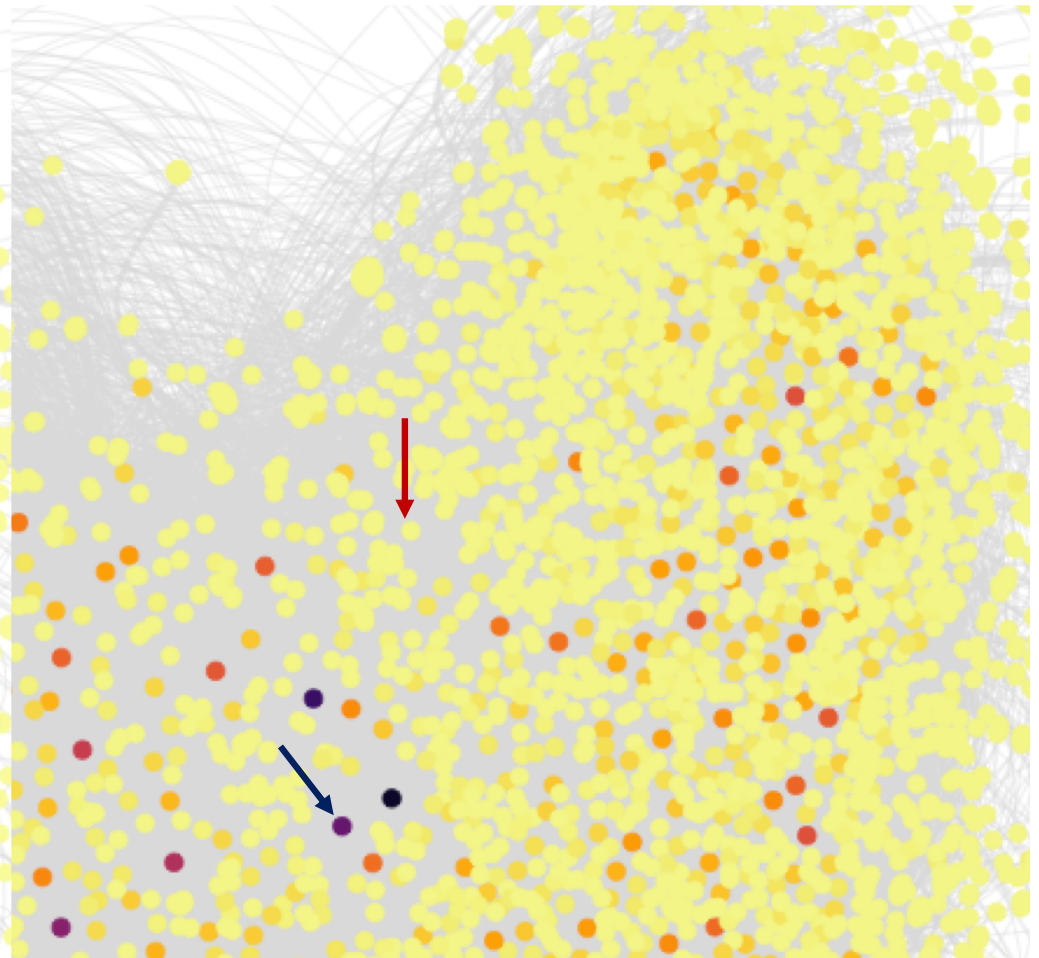
Wiki-vote network example

PageRank centrality

Authorities

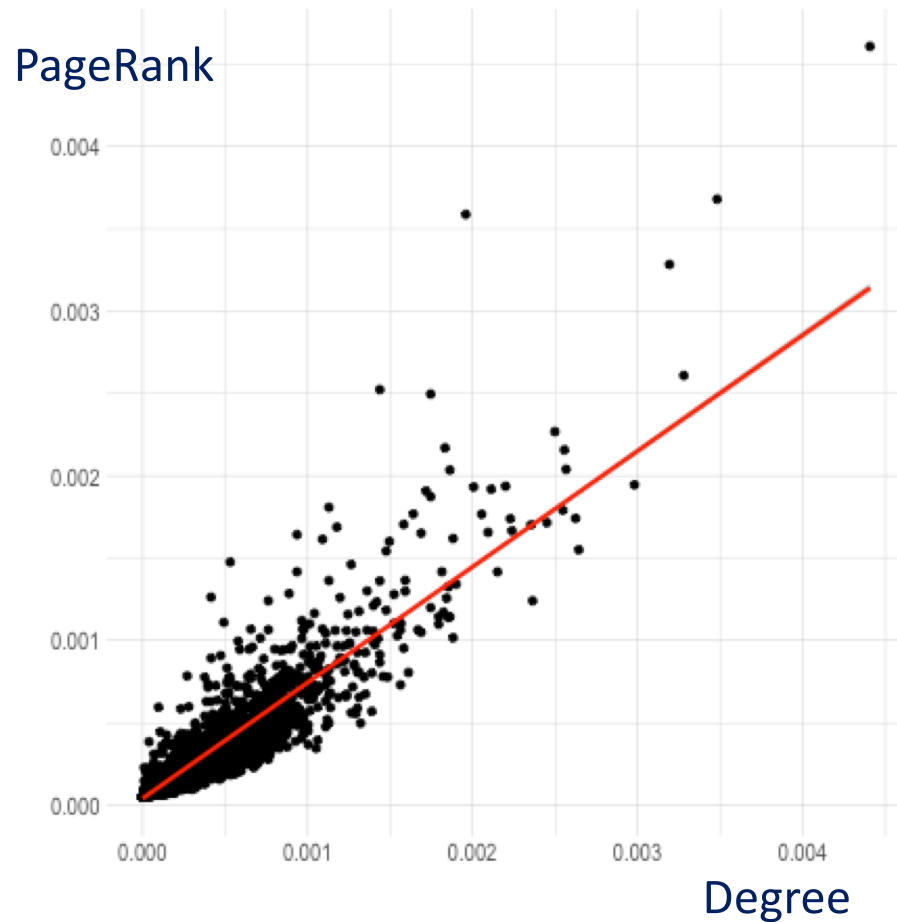


Hubs

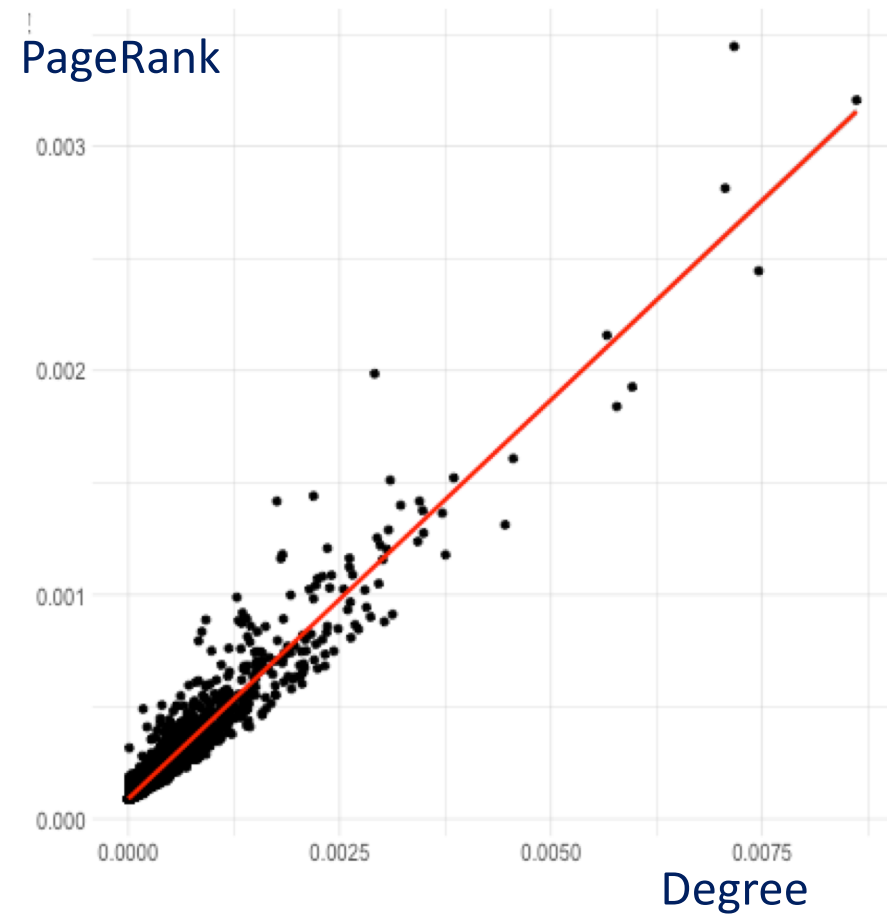


Degree vs PageRank

Authorities

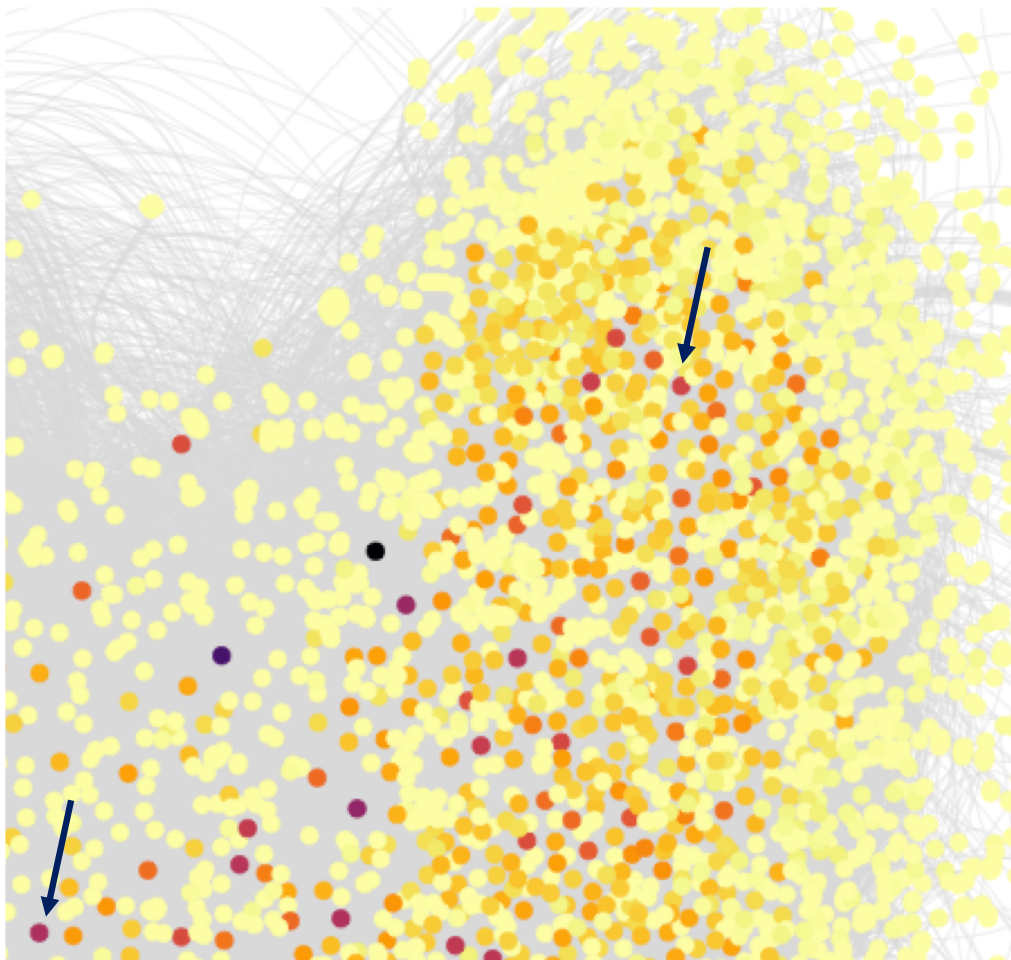


Hubs

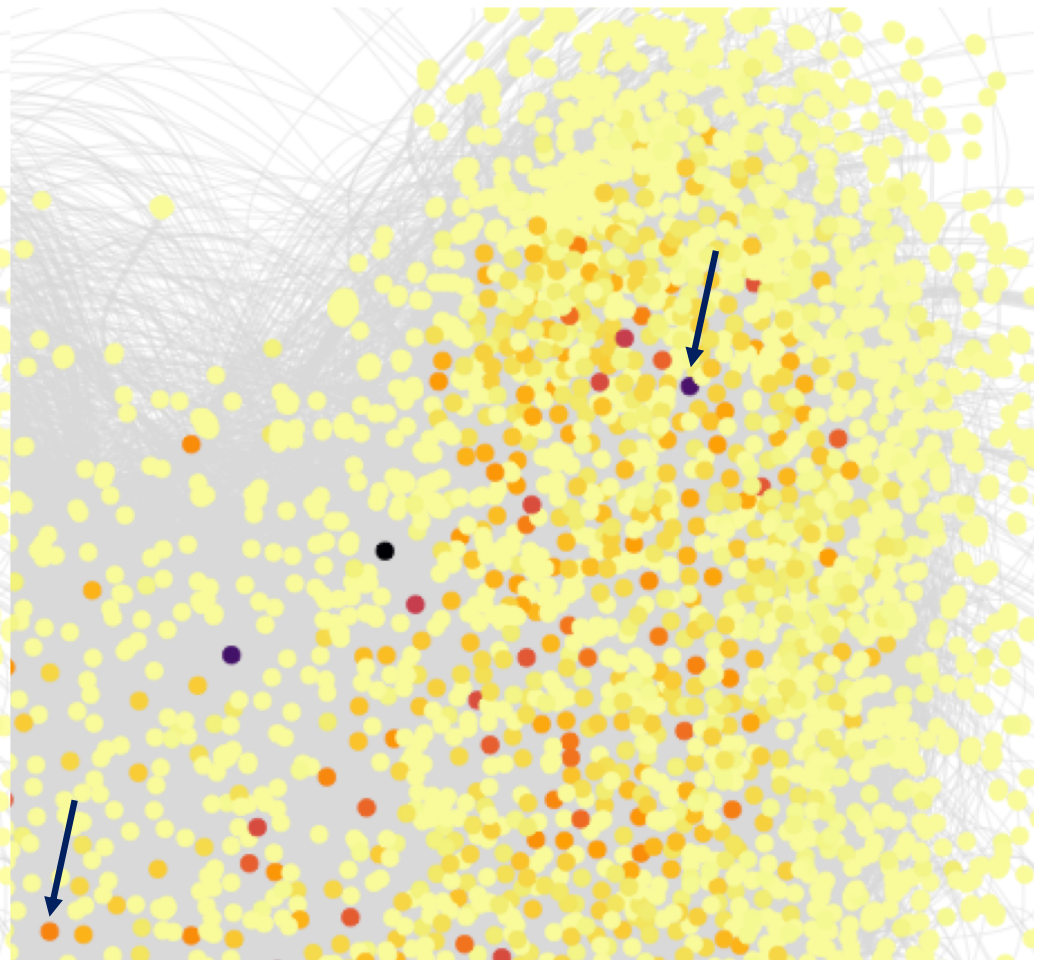


Authorities

Degree

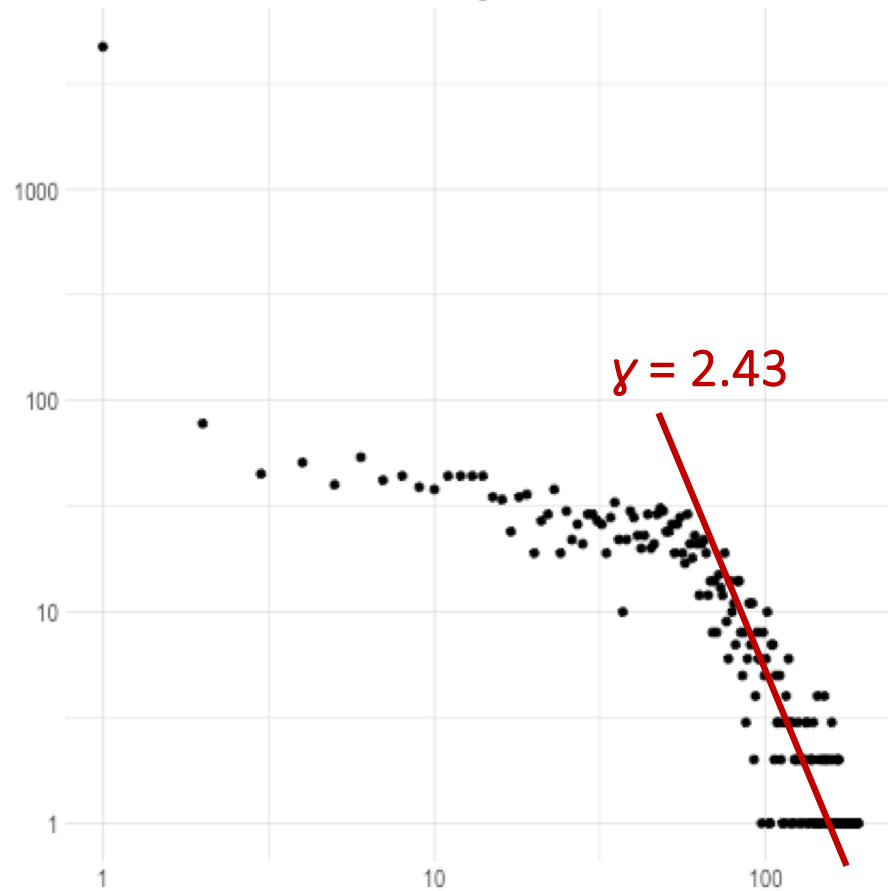


PageRank

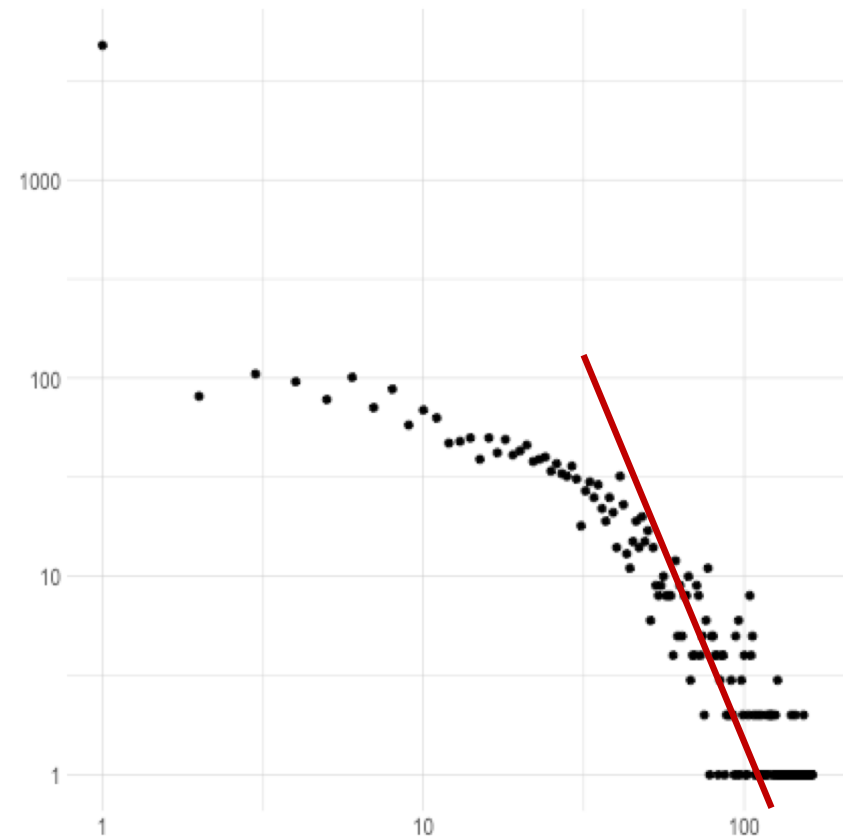


Authorities

Degree

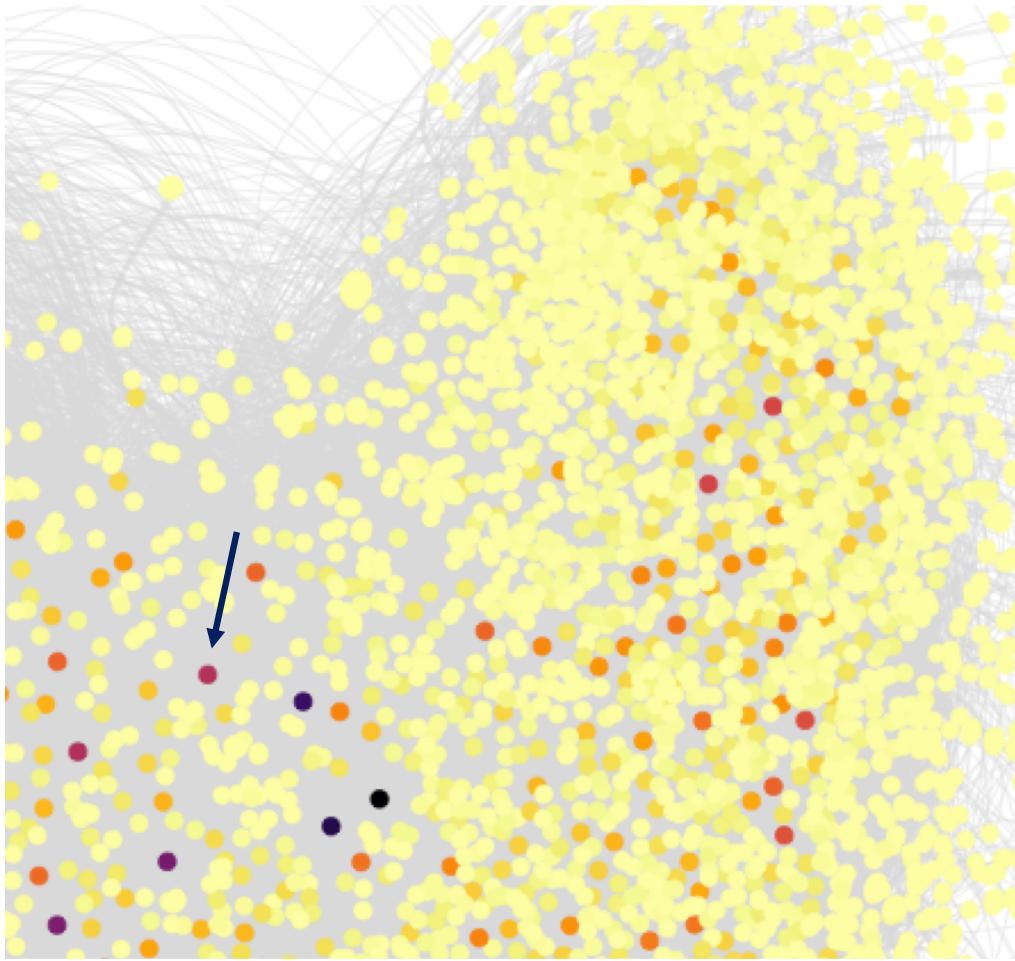


PageRank

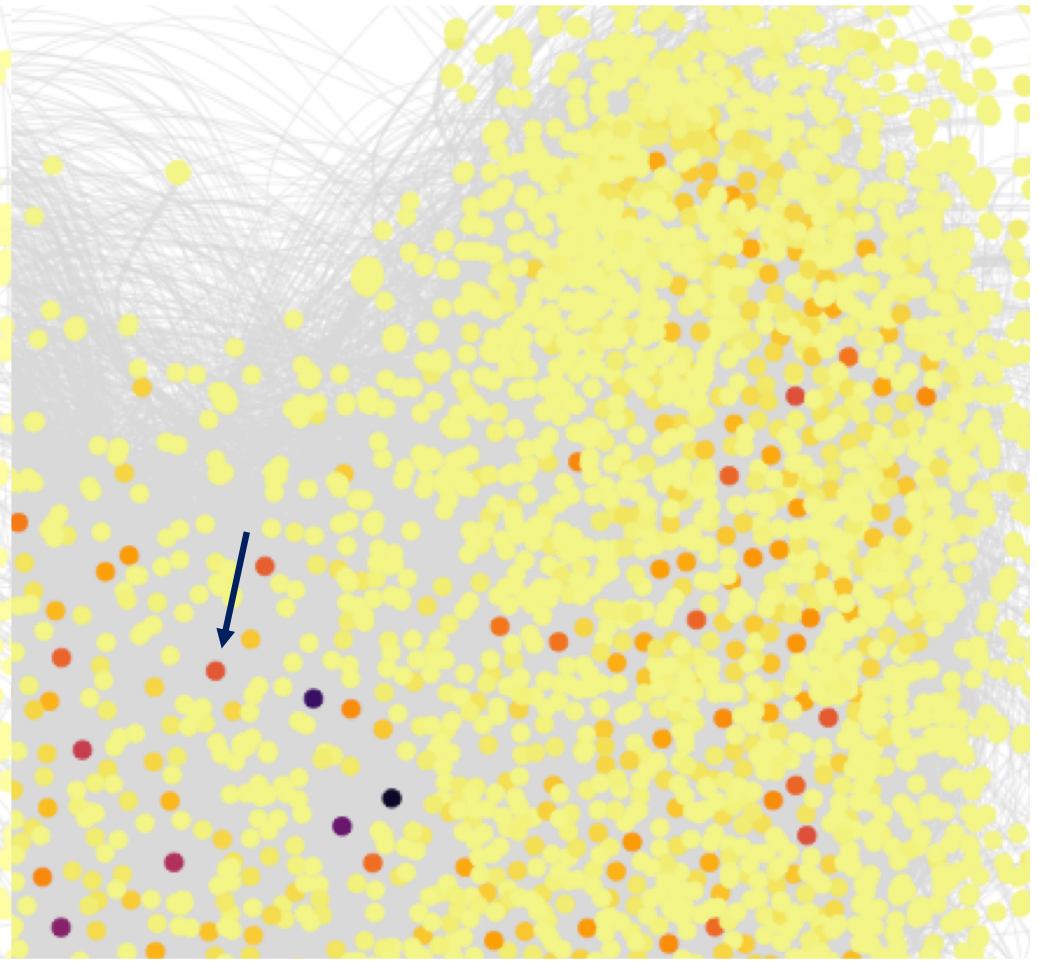


Hubs

Degree

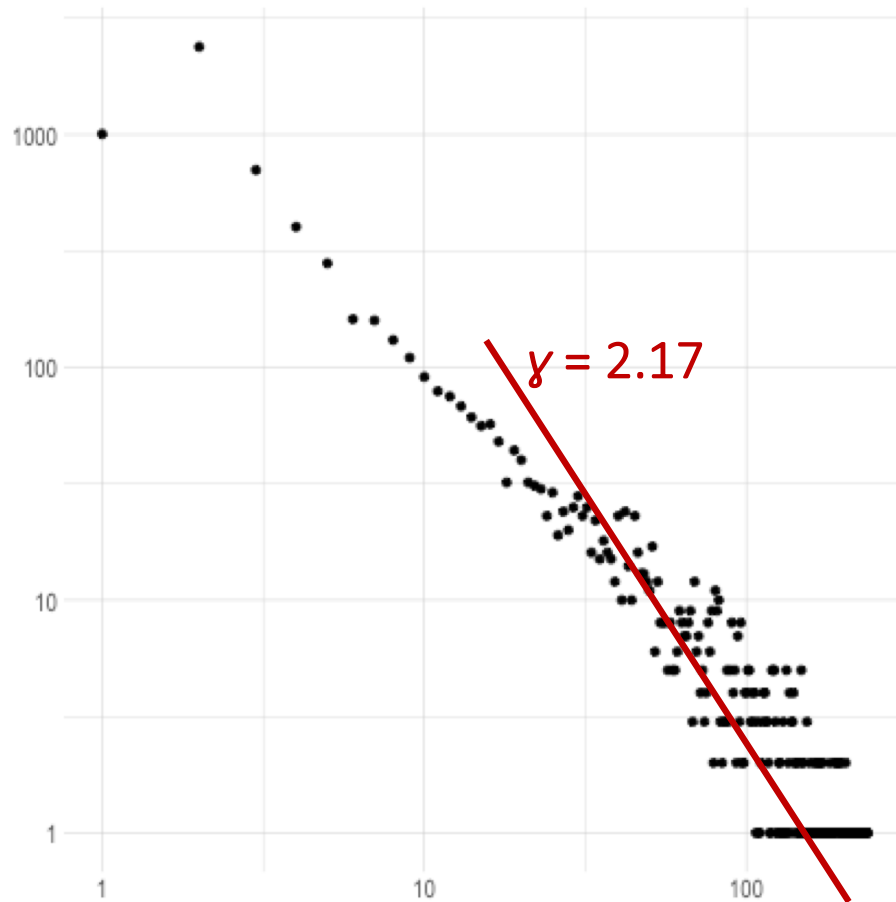


PageRank

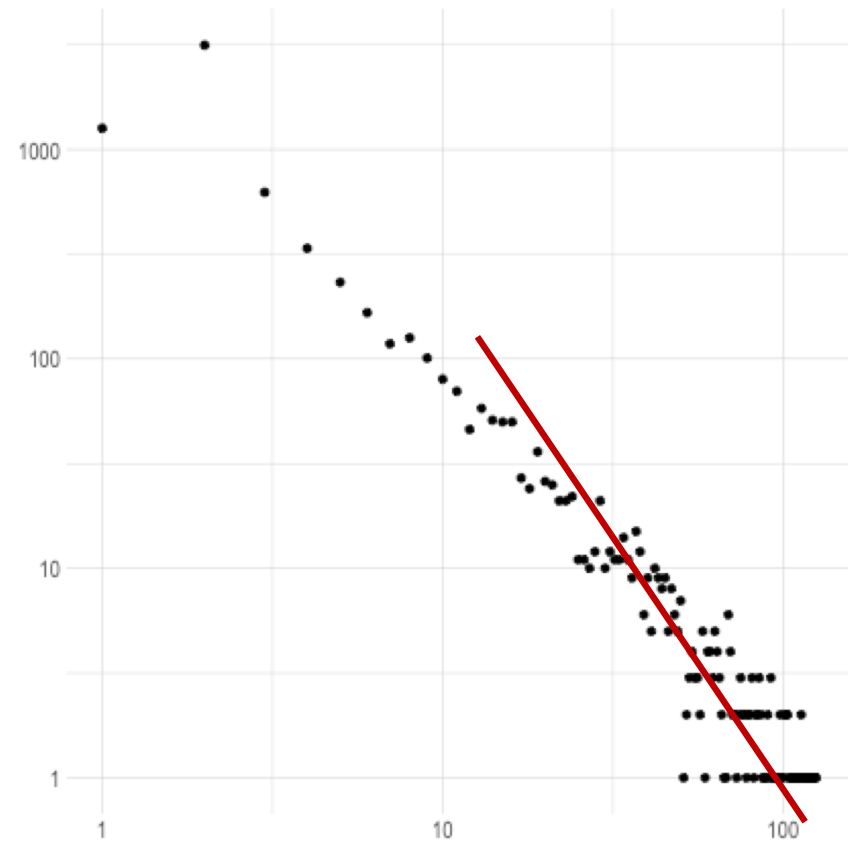


Hubs

Degree



PageRank




Lessons learned

- ❑ Two different approaches
- ❑ Based on simple linear algebra concepts
- ❑ Scalable
- ❑ Implementable via simple message exchange algorithms

HITS

$$h = c M h$$
$$a = c_a A h$$

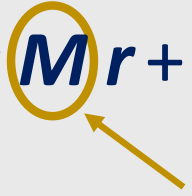
$A^T A$



PageRank

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

$A \text{diag}^{-1}(d)$



Readings

HITS

- Kleinberg, “Authoritative sources in a hyperlinked environment,” 1999

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

PageRank

- 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>

Power iteration

- Wikipedia, “Power iteration”



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