

# A Weighted Correlation Index for Rankings with Ties

Sebastiano Vigna  
Università degli Studi di Milano

# The problem

- Understanding the correlation between different rankings
- Why do correlation measurements between centrality measures are so flaky?
- Taking care of ties is essential (indegree)
- Rank differences between important elements *should be more relevant*
- Large-scale target (whole graphs), not small sets of results

# As a Motivation

- Indegree
- Bavelas's Closeness
- Harmonic Centrality
- PageRank
- Katz

# Geometric Centralities

- Bavelas's closeness (1948):  $\frac{1}{\sum_y d(y, x)}$
- Harmonic centrality (1965):  $\sum_{y \neq x} \frac{1}{d(y, x)}$
- We're just moving from a denormalized, reciprocated arithmetic mean to a denormalized, reciprocated harmonic mean

# Spectral Centralities

- Katz (1951):  $\mathbf{1} \sum_{k \geq 0} \alpha^k G^k$
- PageRank:  $\mathbf{1}/n \sum_{k \geq 0} \alpha^k \bar{G}^k$

Indegree	PageRank	Katz	Harmonic	Closeness
United States	United States	United States	United States	<b>Kharqan Rural District</b>
List of sovereign states	Animal	List of sovereign states	United Kingdom	<b>Talageh-ye Sofla</b>
Animal	List of sovereign states	United Kingdom	World War II	<b>Talageh-ye Olya</b>
England	France	France	France	<b>Greatest Remix Hits (Whigfield album)</b>
France	Germany	Animal	Germany	<b>Suzhou HSR New Town</b>
Association football	Association football	World War II	Association football	<b>Suzhou Lakeside New City</b>
United Kingdom	England	England	English language	<b>Mepirodipine</b>
Germany	India	Association football	China	<b>List of MPs ... M–N</b>
Canada	United Kingdom	Germany	Canada	<b>List of MPs ... O–R</b>
World War II	Canada	Canada	India	<b>List of MPs ... S–T</b>
India	Arthropod	India	<b>Latin</b>	<b>List of MPs ... U–Z</b>
Australia	Insect	Australia	World War I	<b>List of MPs ... J–L</b>
London	World War II	London	England	<b>List of MPs ... C</b>
Japan	Japan	Italy	Italy	<b>List of MPs ... F–I</b>
Italy	Australia	Japan	<b>Russia</b>	<b>List of MPs ... A–B</b>
Arthropod	Village	New York City	<b>Europe</b>	<b>List of MPs ... D–E</b>
Insect	Italy	English language	Australia	<b>Esmaili-ye Sofla</b>
New York City	Poland	China	<b>European Union</b>	<b>Esmaili-ye Olya</b>
English language	English language	Poland	<b>Catholic Church</b>	<b>Levels of organization (ecology)</b>
Village	<b>Nationa Reg. of Hist. Places</b>	World War I	London	<b>Jacques Moeschal (architect)</b>

Table 1: Top 20 pages of the English version of Wikipedia following five different centrality measures.

	Ind.	PR	Katz	Harm.	Cl.		Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.75	0.90	0.62	0.55	Indegree	1	0.31	0.63	0.24	0.06
PageRank	0.75	1	0.75	0.61	0.56	PageRank	0.31	1	0.27	0.10	0.10
Katz	0.90	0.75	1	0.70	0.62	Katz	0.63	0.27	1	0.50	0.20
Harmonic	0.62	0.61	0.70	1	0.92	Harmonic	0.24	0.10	0.50	1	0.65
Closeness	0.55	0.56	0.62	0.92	1	Closeness	0.06	0.10	0.20	0.65	1



Indegree	PageRank	Katz	Harmonic	Closeness
Carl Linnaeus	Carl Linnaeus	Carl Linnaeus	Aristotle	<b>Noël Bernard (botanist)</b>
Aristotle	Aristotle	Aristotle	Albert Einstein	<b>Charles Coquelin</b>
Thomas Jefferson	Thomas Jefferson	Thomas Jefferson	Thomas Jefferson	<b>Markku Kivinen</b>
Margaret Thatcher	Charles Darwin	Albert Einstein	Charles Darwin	<b>Angiolo Maria Colomboni</b>
Plato	Plato	Charles Darwin	Thomas Edison	<b>Om Prakash (historian)</b>
Charles Darwin	Albert Einstein	Karl Marx	<b>Alexander Graham Bell</b>	<b>Michel Mandjes</b>
Karl Marx	Karl Marx	Plato	<b>Nikola Tesla</b>	<b>Kees Posthumus</b>
Albert Einstein	Pliny the Elder	Margaret Thatcher	<b>William James</b>	<b>F. Wolfgang Schnell</b>
Vladimir Lenin	Vladimir Lenin	Vladimir Lenin	Isaac Newton	<b>Christof Ebert</b>
Sigmund Freud	Johann Wolfgang von Goethe	Isaac Newton	Karl Marx	<b>Reese Prosser</b>
J. R. R. Tolkien	Margaret Thatcher	Ptolemy	<b>Charles Sanders Peirce</b>	<b>David Tulloch</b>
Johann Wolfgang von Goethe	Ptolemy	Johann Wolfgang von Goethe	Noam Chomsky	<b>Kim Hawtrey</b>
<b>Spider-Man</b>	Sigmund Freud	Pliny the Elder	<b>Enrico Fermi</b>	<b>Patrick J. Miller</b>
Pliny the Elder	Isaac Newton	Benjamin Franklin	Ptolemy	<b>Mikel King</b>
Benjamin Franklin	Benjamin Franklin	J. R. R. Tolkien	<b>John Dewey</b>	<b>Albert Perry Brigham</b>
Leonardo da Vinci	J. R. R. Tolkien	Thomas Edison	Johann Wolfgang von Goethe	<b>Gordon Wagner (economist)</b>
Isaac Newton	Immanuel Kant	Sigmund Freud	<b>Bertrand Russell</b>	<b>George Henry Chase</b>
Ptolemy	Leonardo da Vinci	Immanuel Kant	Plato	<b>Charles C. Horn</b>
Immanuel Kant	<b>Pierre André Latreille</b>	Leonardo da Vinci	<b>John von Neumann</b>	<b>Paul Goldstene</b>
<b>George Bernard Shaw</b>	Thomas Edison	Noam Chomsky	Vladimir Lenin	<b>Robert Stanton Avery</b>

Indegree	PageRank	Katz	Harmonic	Closeness
Martini (cocktail)	Martini (cocktail)	Irish coffee	Irish coffee	<b>Magie Noir</b>
Piña colada	Caipirinha	Caipirinha	Caipirinha	<b>Batini (drink)</b>
Mojito	Mojito	Martini (cocktail)	Kir (cocktail)	<b>Scorpion bowl</b>
Caipirinha	Piña colada	Piña colada	Martini (cocktail)	<b>Poinsettia (cocktail)</b>
Cuba Libre	Irish coffee	Kir (cocktail)	Piña colada	Irish coffee
Irish coffee	Kir (cocktail)	Mojito	Mojito	Caipirinha
Singapore Sling	Cosmopolitan (cocktail)	Mai Tai	Beer cocktail	Kir (cocktail)
Manhattan (cocktail)	Manhattan (cocktail)	Cuba Libre	Shaken, not stirred	Martini (cocktail)
Windle (sidecar)	IBA Official Cocktail	Singapore Sling	Pisco Sour	Piña colada
Cosmopolitan (cocktail)	Beer cocktail	Long Island Iced Tea	Mai Tai	Mojito
Mai Tai	Mai Tai	Shaken, not stirred	Spritz (alcoholic beverage)	Beer cocktail
IBA Official Cocktail	Singapore Sling	Beer cocktail	Long Island Iced Tea	Shaken, not stirred
Kir (cocktail)	Cuba Libre	Manhattan (cocktail)	Sazerac	Mai Tai
Shaken, not stirred	<b>Tom Collins</b>	Cosmopolitan (cocktail)	Fizz (cocktail)	Spritz (alcoholic beverage)
Beer cocktail	Long Island Iced Tea	Windle (sidecar)	Flaming beverage	Pisco Sour
Pisco Sour	Sour (cocktail)	Pisco Sour	Cuba Libre	Long Island Iced Tea
Long Island Iced Tea	Shaken, not stirred	White Russian (cocktail)	Wine cocktail	Sazerac
Sour (cocktail)	<b>Negroni</b>	IBA Official Cocktail	Singapore Sling	Flaming beverage
White Russian (cocktail)	Flaming beverage	Moscow mule	Moscow mule	Fizz (cocktail)
Vesper (cocktail)	<b>Lillet</b>	Vesper (cocktail)	White Russian (cocktail)	Wine cocktail

# Kendall's $\tau$ 1938

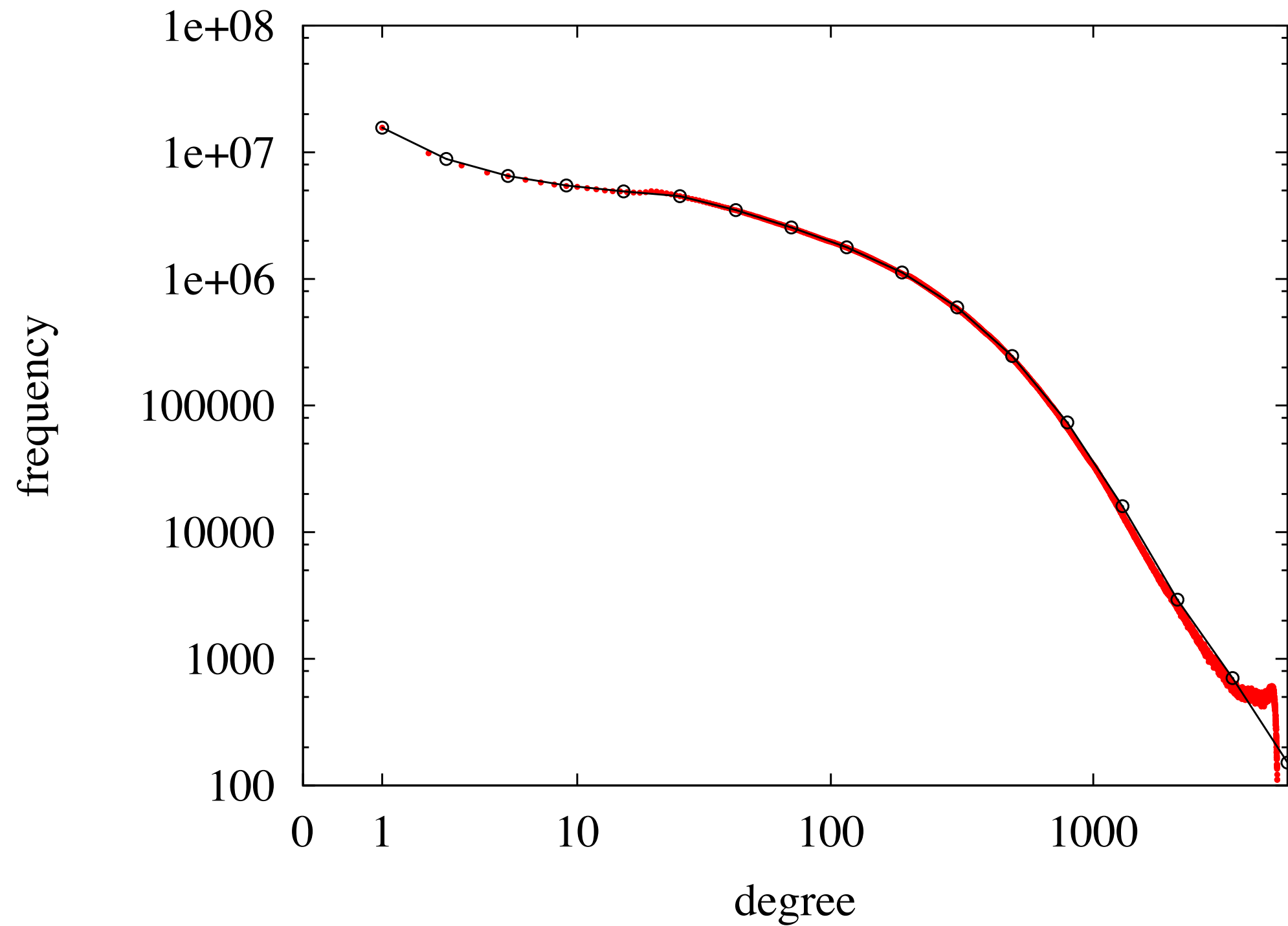
- Scores  $\mathbf{r}, \mathbf{s}$  (distinct)
- Concordances: pairs  $(i, j), i < j$ , such that the scores for  $i$  and  $j$  in  $\mathbf{r}$  and  $\mathbf{s}$  are in the same order
- $\tau$ : Concordances minus discordances divided by concordances plus discordances (i.e., the number of ordered pairs)
- Note: if you skip ties on both sides you get Goodman–Kruskal's  $\gamma$



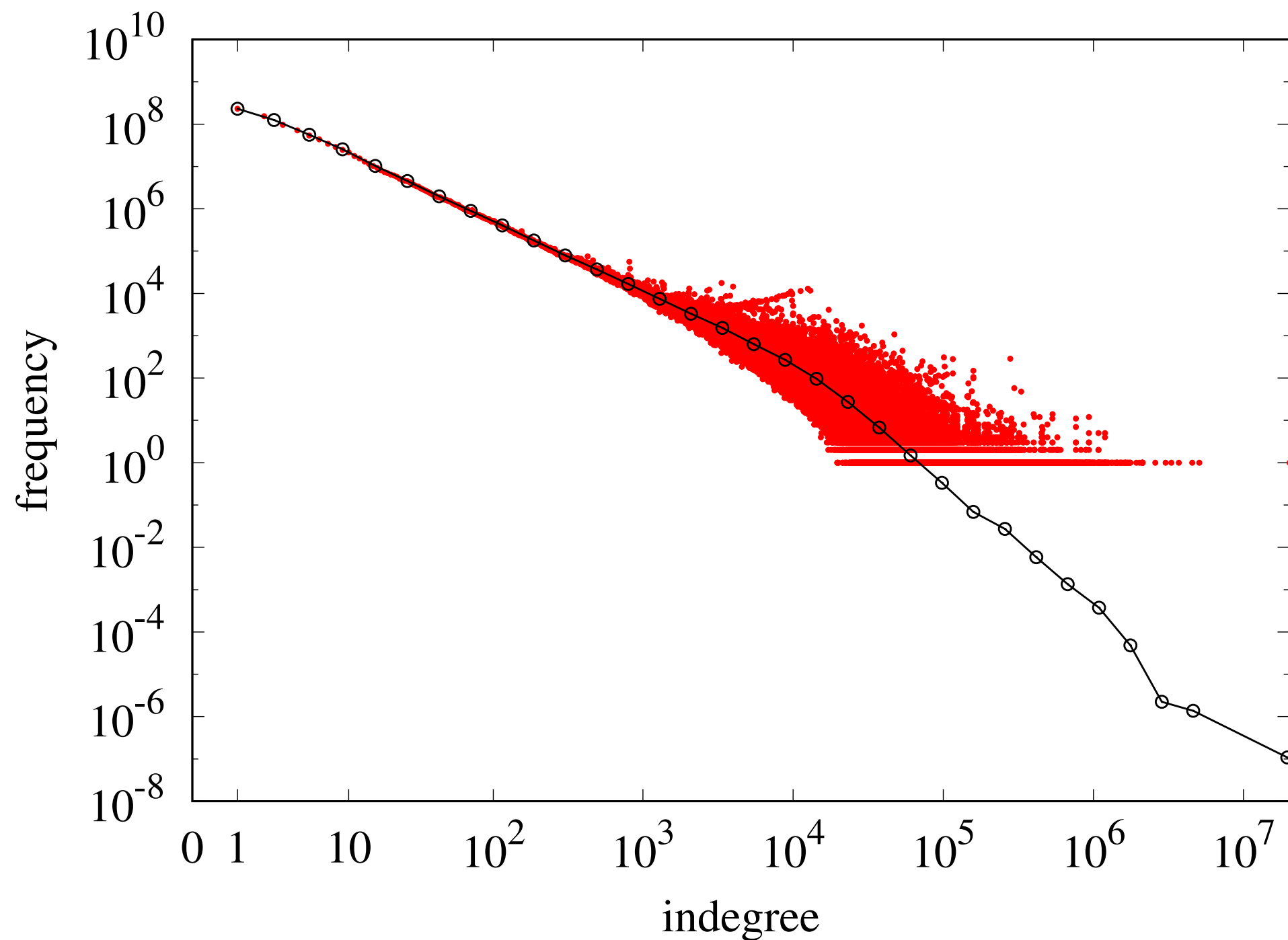
# Ties Are Important

- >99% nodes in a typical social/web graph are in a degree tie
- Ties cannot be solved by random assignment
- $\langle 0, 0, 0, \dots, 1, 1, 1, \dots \rangle$  and  $\langle 1, 1, 1, \dots, 2, 2, 2, \dots \rangle$  give correlation  $\approx 0.5$
- When you compare network rankings, you're almost always working in the 0.7-to-1 region

# Facebook (2011)



• .eu (2015,  $10^9$  pages)



# Kendall's $\tau$ 1945

- Starts from Daniels (1944): every correlation is a cosine similarity

$$\langle \mathbf{r}, \mathbf{s} \rangle := \sum_{i < j} \text{sgn}(r_i - r_j) \text{sgn}(s_i - s_j)$$

- “Norm”

$$\|\mathbf{r}\| := \sqrt{\langle \mathbf{r}, \mathbf{r} \rangle}$$

- Cosine similarity!

$$\tau(\mathbf{r}, \mathbf{s}) := \frac{\langle \mathbf{r}, \mathbf{s} \rangle}{\|\mathbf{r}\| \cdot \|\mathbf{s}\|}$$

# Now

$$\langle \mathbf{r}, \mathbf{s} \rangle_w := \sum_{i < j} \operatorname{sgn}(r_i - r_j) \operatorname{sgn}(s_i - s_j) w(i, j)$$

$$|\langle \mathbf{r}, \mathbf{s} \rangle_w| \leq \|\mathbf{r}\|_w \|\mathbf{s}\|_w$$

$$\tau_w(\mathbf{r}, \mathbf{s}) := \frac{\langle \mathbf{r}, \mathbf{s} \rangle_w}{\|\mathbf{r}\|_w \cdot \|\mathbf{s}\|_w}$$

# Related

- Shieh [St. & Pr. L. 1998]: weight the pair  $(i, j)$  with  $w_{ij}$
- No ties!
- Yilmaz, Aslam & Robertson [SIGIR 2008]:  $w_{ij} := 1/j$  (motivated by probability)
- Kumar & Vassilvitskii [WWW 2010]:  $w_{ij}$  depending on similarity, ties broken randomly
- Iman & Conover [Technometrics 1987]: Spearman's correlation between *Savage* scores  $H_n - H_{i-1}$
- Farnoud [2012]: weight adjacent transpositions  $(i \ i + 1)$  and compute minimum weight

# Decoupling Rank and Weight

$$\langle \mathbf{r}, \mathbf{s} \rangle_{\rho, w} := \sum_{i < j} \text{sgn}(r_i - r_j) \text{sgn}(s_i - s_j) w(\rho(i), \rho(j))$$

$$\tau_{\bullet, w}(\mathbf{r}, \mathbf{s}) := \frac{\tau_{\rho_{\mathbf{r}, \mathbf{s}}, w}(\mathbf{r}, \mathbf{s}) + \tau_{\rho_{\mathbf{s}, \mathbf{r}}, w}(\mathbf{r}, \mathbf{s})}{2}$$

- $\rho_{\mathbf{r}, \mathbf{s}}$  is the ranking induced by  $\mathbf{r}$  and  $\mathbf{s}$  in lexicographical order
- viceversa for  $\rho_{\mathbf{s}, \mathbf{r}}$

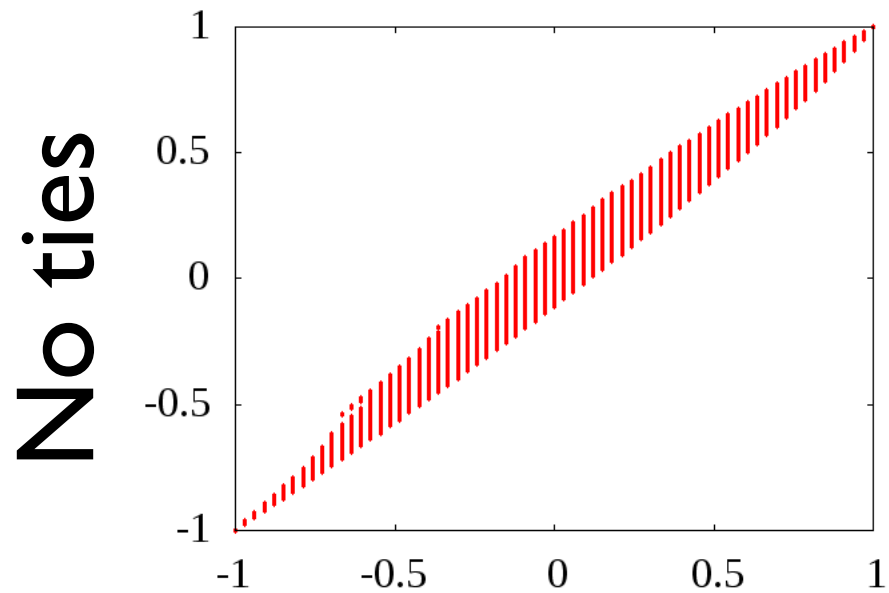


# Computable Quickly

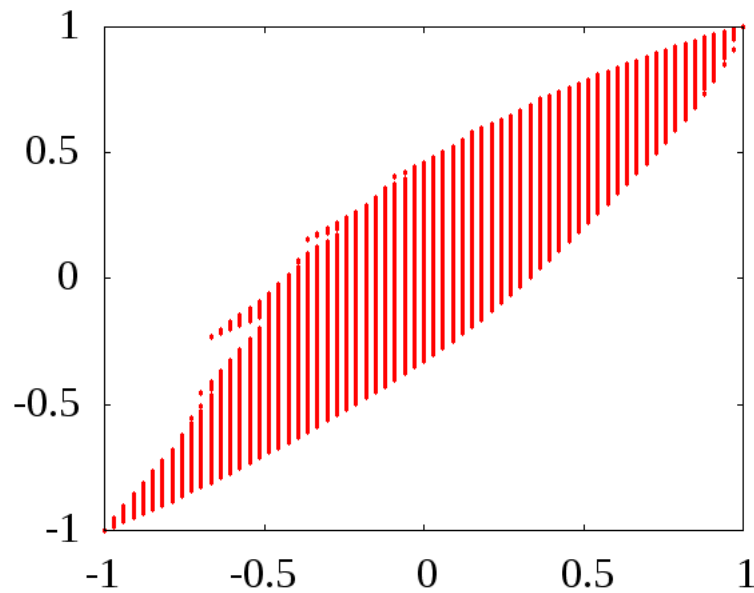
- $O(n \log n)$  variant of Knight's algorithm (highly parallelizable, distributable—it's a MergeSort)
- Works for any scheme  $w(i, j) := f(i) \odot g(j)$  with suitable operation  $\odot$  (e.g., addition, multiplication)
- We suggest *additive hyperbolic* weighting, weighting  $(i, j)$  by  $1/(i+1) + 1/(j+1)$ :  $\tau_h$

# Correlation with Kendall's $\tau$

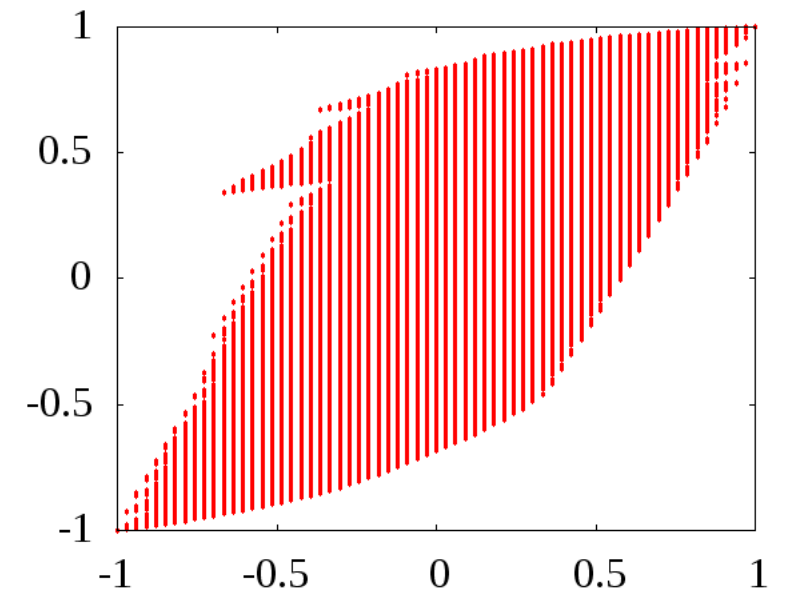
## Logarithmic



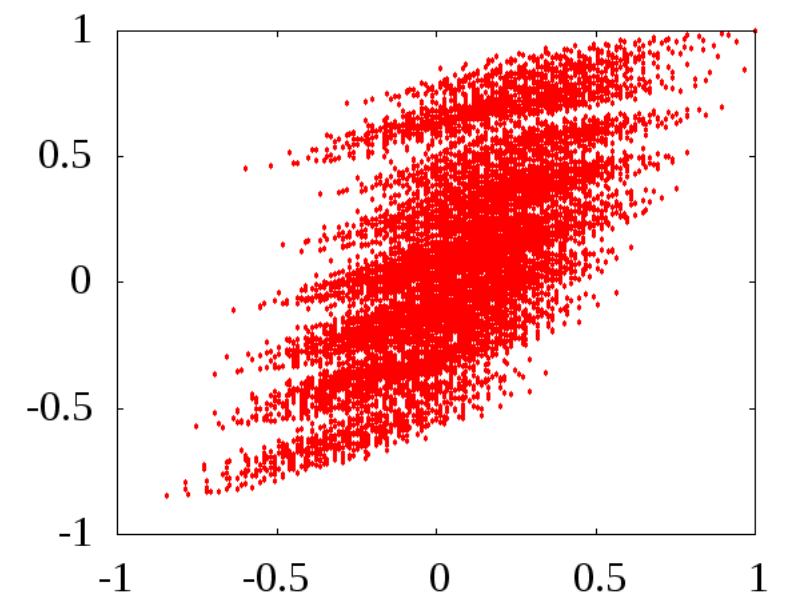
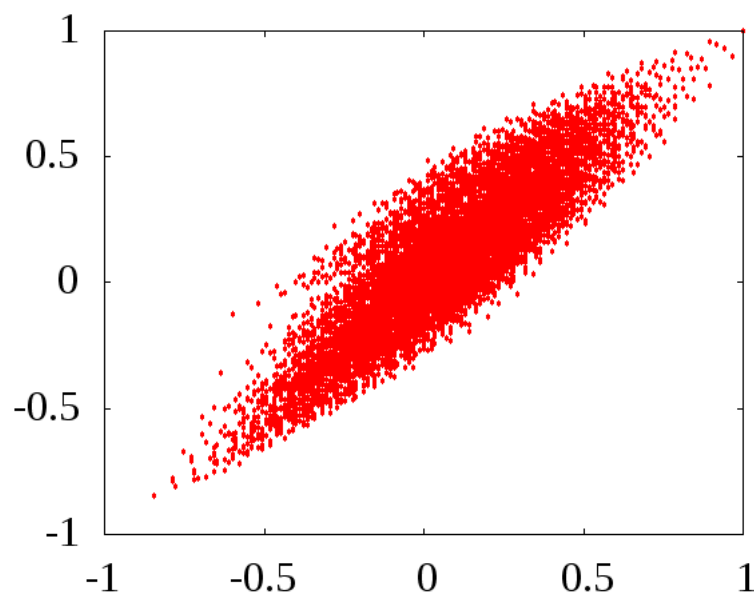
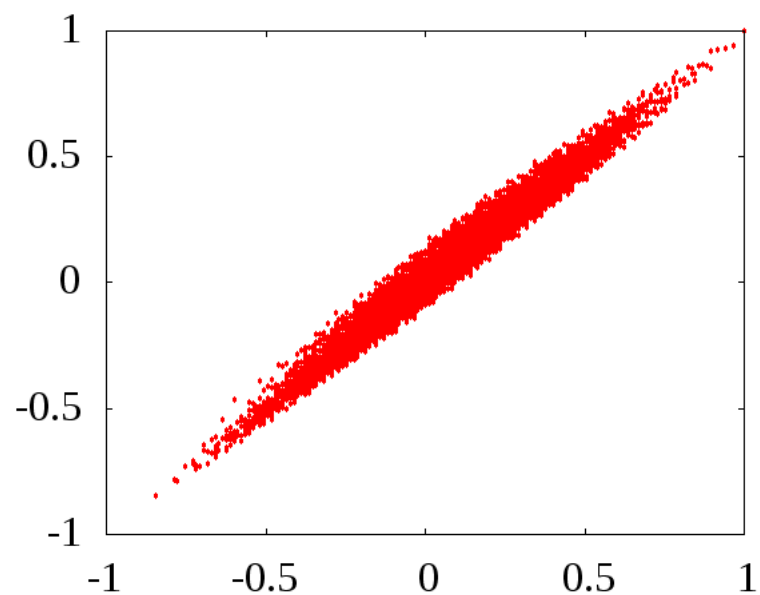
## Hyperbolic



## Quadratic



Ties



# Wikipedia

$\tau$

	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.75	0.90	0.62	0.55
PageRank	0.75	1	0.75	0.61	0.56
Katz	0.90	0.75	1	0.70	0.62
Harmonic	0.62	0.61	0.70	1	0.92
Closeness	0.55	0.56	0.62	0.92	1

$\tau_h$

	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.95	0.98	0.90	0.27
PageRank	0.95	1	0.96	0.92	0.65
Katz	0.98	0.96	1	0.93	0.26
Harmonic	0.90	0.92	0.93	1	0.28
Closeness	0.27	0.65	0.26	0.28	1

Table 6:  $\tau_h$  on Wikipedia.

# Hollywood co-starship

Indegree	PageRank	Katz	Harmonic	Closeness
Shatner, William	Jeremy, Ron	Shatner, William	Sheen, Martin	<b>Östlund, Claes Göran</b>
Flowers, Bess	Hitler, Adolf	Sheen, Martin	Clooney, George	<b>Östlund, Catarina</b>
Sheen, Martin	<b>Kaufman, Lloyd</b>	Hanks, Tom	Jackson, Samuel L.	<b>von Preußen, Oskar Prinz</b>
Reagan, Ronald (I)	Bush, George W.	Williams, Robin (I)	Hopper, Dennis	<b>von Preußen, Georg Friedrich</b>
Clooney, George	Reagan, Ronald (I)	Clooney, George	Hanks, Tom	<b>von Mannstein, Robert Grund</b>
Jackson, Samuel L.	Clinton, Bill (I)	Reagan, Ronald (I)	Stone, Sharon (I)	<b>von Mannstein, Concha</b>
Williams, Robin (I)	Sheen, Martin	Willis, Bruce	Brosnan, Pierce	<b>von der Busken, Mart</b>
Hanks, Tom	<b>Rochon, Debbie</b>	Jackson, Samuel L.	Hitler, Adolf	<b>van der Putten, Thea</b>
Jeremy, Ron	<b>Kennedy, John F.</b>	Stone, Sharon (I)	<b>McDowell, Malcolm</b>	<b>de la Bruheze, Joel Albert</b>
Hitler, Adolf	Hopper, Dennis	Freeman, Morgan (I)	Williams, Robin (I)	<b>de la Bruheze, Emile</b>
Willis, Bruce	<b>Nixon, Richard</b>	Flowers, Bess	<b>De Niro, Robert</b>	<b>te Riele, Marloes</b>
Clinton, Bill (I)	<b>Estevez, Joe</b>	Brosnan, Pierce	Willis, Bruce	<b>de Reijer, Eric</b>
Freeman, Morgan (I)	Shatner, William	Douglas, Michael (I)	<b>Hopkins, Anthony</b>	<b>des Bouvrie, Jan</b>
Hopper, Dennis	Jackson, Samuel L.	Madonna (I)	Madonna (I)	<b>de Klijn, Judith</b>
Stone, Sharon (I)	<b>Stewart, Jon (I)</b>	Travolta, John	<b>Lee, Christopher (I)</b>	<b>de Freitas, Luís (II)</b>
Madonna (I)	<b>Carradine, David (I)</b>	Hopper, Dennis	Douglas, Michael (I)	<b>de Freitas, Luís (I)</b>
Bush, George W.	Asner, Edward	Ford, Harrison (I)	<b>Sutherland, Donald (I)</b>	<b>Zuu, Winnie Otondi</b>
<b>Harris, Sam (II)</b>	<b>Zirnkilton, Steven</b>	Asner, Edward	Freeman, Morgan (I)	<b>Zuu, Emmanuel Dahngbay</b>
Brosnan, Pierce	<b>Colbert, Stephen</b>	<b>MacLaine, Shirley</b>	<b>Stallone, Sylvester</b>	<b>Zilbersmith, Carla</b>
Travolta, John	<b>Madsen, Michael (I)</b>	Clinton, Bill (I)	Ford, Harrison (I)	<b>Zilber, Mac</b>

Table 10: Top 20 pages of the Hollywood co-starship graph.

# Hollywood co-starship

$\tau$

	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.42	0.93	0.55	0.43
PageRank	0.42	1	0.36	0.10	0.18
Katz	0.93	0.36	1	0.61	0.49
Harmonic	0.55	0.10	0.61	1	0.86
Closeness	0.43	0.18	0.49	0.86	1

$\tau_h$

	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.90	0.98	0.91	0.10
PageRank	0.90	1	0.88	0.81	0.64
Katz	0.98	0.88	1	0.92	0.11
Harmonic	0.91	0.81	0.92	1	0.18
Closeness	0.10	0.64	0.11	0.18	1



# Common Crawl Hosts

Indegree	PageRank	Katz	Harmonic	Closeness
wordpress.org	gmpg.org	wordpress.org	youtube.com	0-p.com
youtube.com	wordpress.org	youtube.com	en.wikipedia.org	0-0-0-0-0-0-0.indahiphop.ru
gmpg.org	youtube.com	gmpg.org	twitter.com	0-0-1.i.tiexue.net
en.wikipedia.org	<b>livejournal.com</b>	en.wikipedia.org	google.com	0-00cigarettes.info
tumblr.com	tumblr.com	tumblr.com	wordpress.org	0-0mos00.hi5.com
twitter.com	en.wikipedia.org	twitter.com	flickr.com	0-0new0-0.hi5.com
google.com	twitter.com	google.com	facebook.com	0-0sunny0-0.hi5.com
flickr.com	<b>networkadvertising.org</b>	flickr.com	<b>apple.com</b>	0-1.i.tiexue.net
rtalabel.org	<b>promodj.com</b>	rtalabel.org	vimeo.com	0-1.sxsy.co
wordpress.com	<b>skriptmail.de</b>	wordpress.com	creativecommons.org	0-2.paparazziwannabe.com
mp3shake.com	<b>parallels.com</b>	mp3shake.com	<b>amazon.com</b>	0-311.cn
w3schools.com	<b>tistory.com</b>	w3schools.com	<b>adobe.com</b>	0-360.rukazan.ru
domains.lycos.com	google.com	creativecommons.org	<b>myspace.com</b>	0-5days.com
staff.tumblr.com	miibeian.gov.cn	staff.tumblr.com	<b>w3.org</b>	0-5days.net
club.tripod.com	phpbb.com	domains.lycos.com	<b>bbc.co.uk</b>	0-5kalibr.pdj.ru
creativecommons.org	<b>blog.fc2.com</b>	club.tripod.com	<b>nytimes.com</b>	0-9-0-4-4-9.promoradio.ru
vimeo.com	<b>tw.yahoo.com</b>	vimeo.com	<b>yahoo.com</b>	0-9-0-9.dbass.ru
miibeian.gov.cn	w3schools.com	miibeian.gov.cn	<b>microsoft.com</b>	0-9-0-9.promodj.ru
facebook.com	wordpress.com	facebook.com	<b>guardian.co.uk</b>	0-9-1125.i.tiexue.net
phpbb.com	domains.lycos.com	phpbb.com	<b>imdb.com</b>	0-9-7-16.software.informer.com

# Common Crawl Hosts

$\tau$

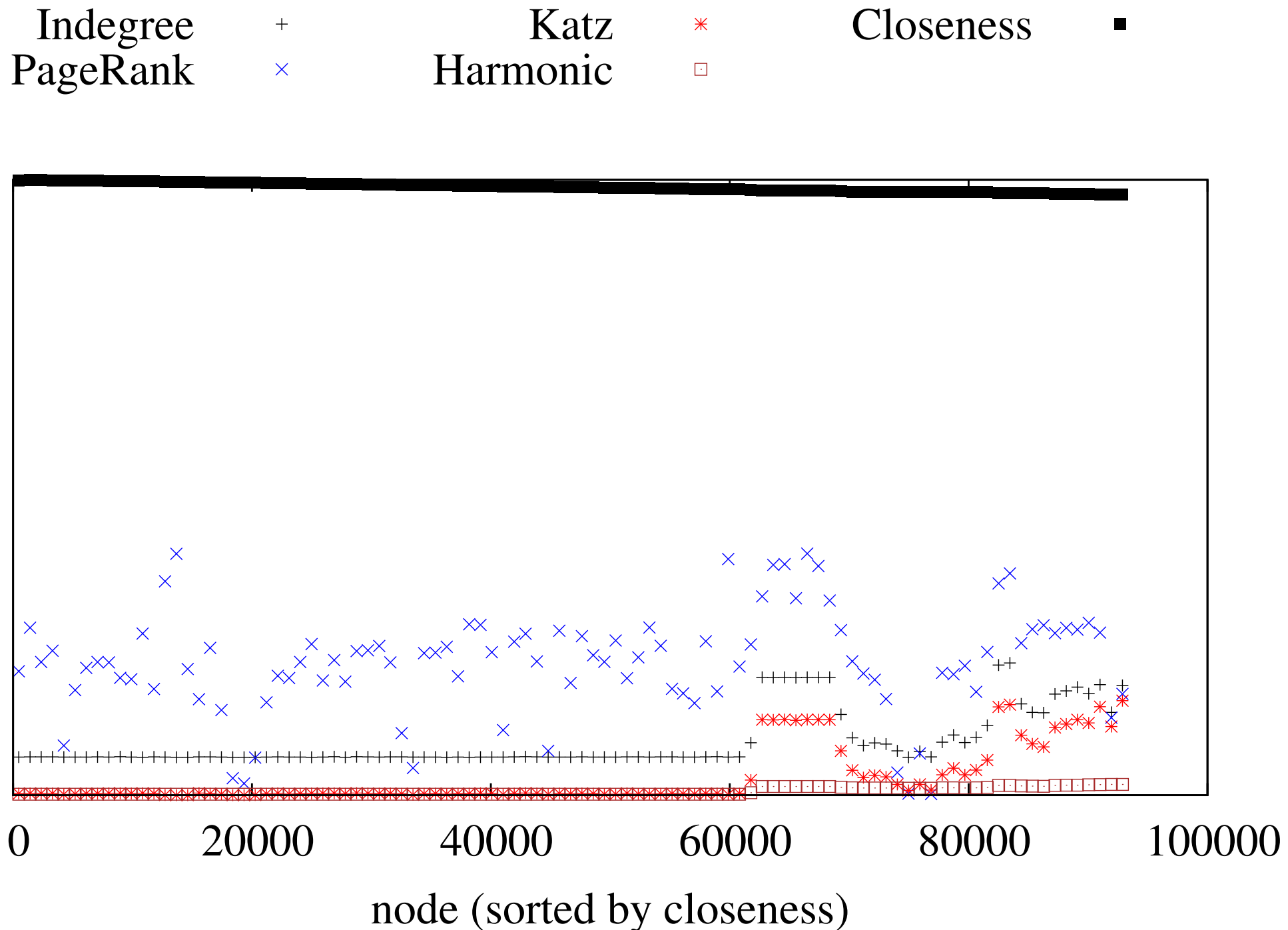
	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.71	0.89	0.61	0.54
PageRank	0.71	1	0.66	0.50	0.50
Katz	0.89	0.66	1	0.69	0.59
Harmonic	0.61	0.50	0.69	1	0.86
Closeness	0.54	0.50	0.59	0.86	1

$\tau_h$

	Ind.	PR	Katz	Harm.	Cl.
Indegree	1	0.91	0.96	0.72	0.20
PageRank	0.91	1	0.90	0.81	0.69
Katz	0.96	0.90	1	0.78	0.15
Harmonic	0.72	0.81	0.78	1	0.35
Closeness	0.20	0.69	0.15	0.35	1

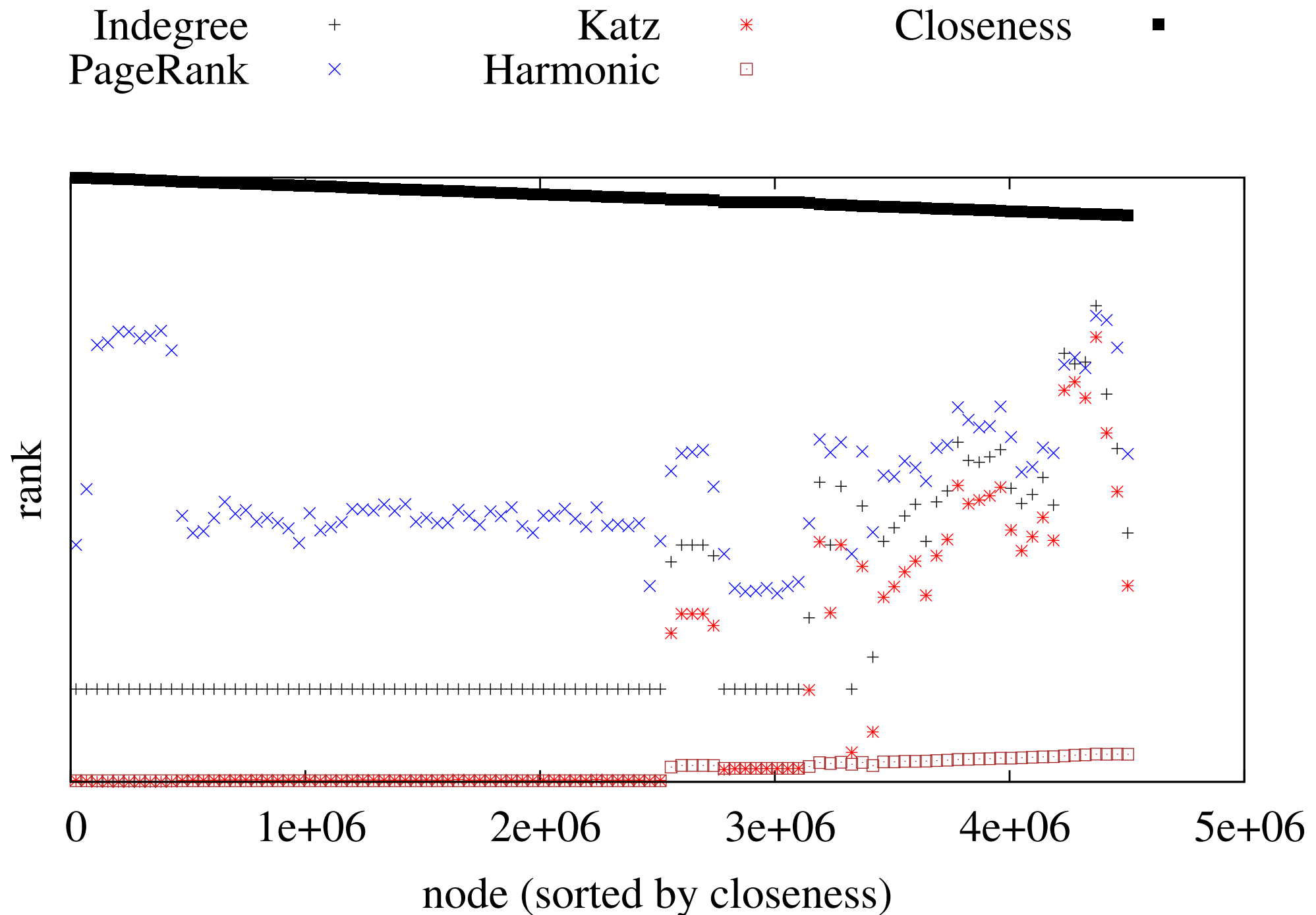


# Why does closeness $T_h$ -correlate with PageRank?



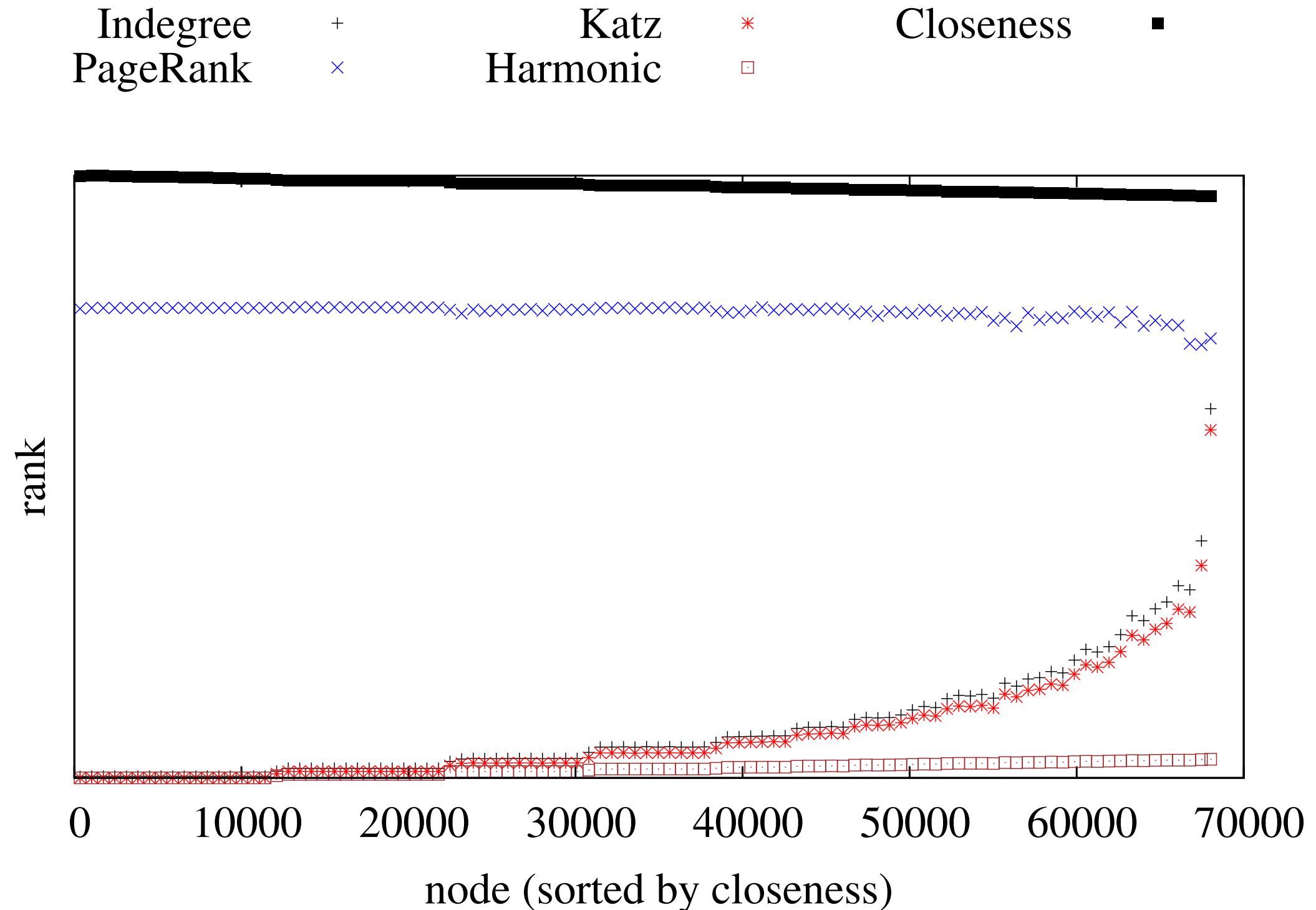
Rank of nodes unreachable from the giant component of Wikipedia

# Why does closeness $T_h$ -correlate with PageRank?



Rank of nodes unreachable from the giant component of Common Crawl

# Why does closeness $T_h$ -correlate with PageRank?



Rank of nodes unreachable from the giant component of Hollywood

# Conclusions

- We believe  $\tau_h$  is a new and valuable tool to understand rankings
- In general, the machinery around  $\tau_w$  can be used to try easily new application-dependent weighted correlation indices on a large scale
- There are obvious elements of arbitrariness, but hyperbolic weighting is at the convergence of several previous proposals
- Implemented as `stats.weightedtau` in SciPy (easy to use!)
- Software, as usual, at <http://law.di.unimi.it/>