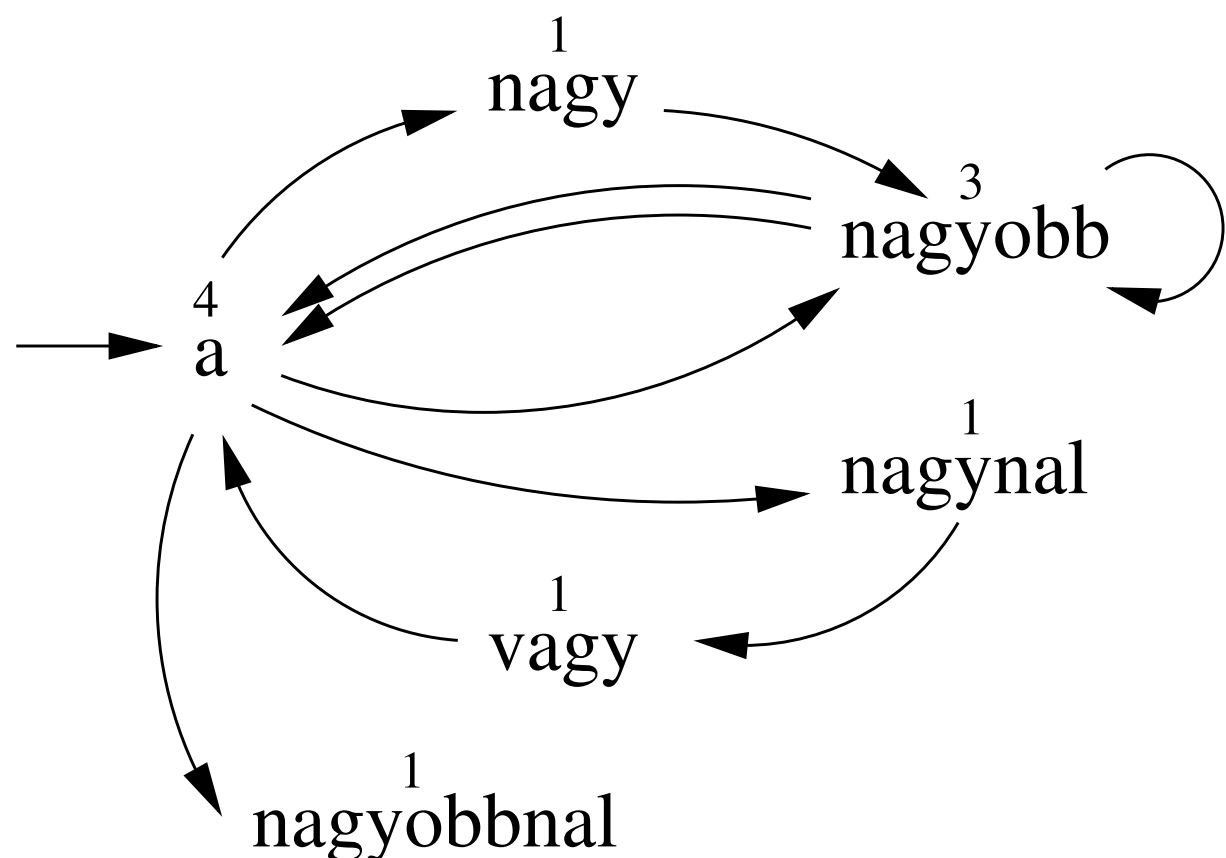


Makrai Márton és Sass Bálint
MTA Nyelvtudományi Intézet
{makrai.marton,sass.balint}@nytud.mta.hu



- szógyakoriságok Zipf (1935)
- skálafüggetlen gráf (Barabási and Albert, 1999)
- most: irányított gráf súlyozott élekkel bigramgyakoriságokból



csúcsok száma	élek száma
0	0
5000	12500
10000	20000
15000	28750
20000	37500
25000	46250
30000	55000
35000	63750
40000	72500
45000	81250

log(cúscok száma)	irányított	irányítatlan
10 ⁰	~7.0	~4.2
10 ¹	~5.9	~4.0
10 ²	~4.1	~3.1

A log-log plot showing the relationship between the number of peaks (log(csúcsok száma) on the x-axis) and transitivity (log(transzitivitás) on the y-axis). The x-axis ranges from 10^0 to 10^2 , and the y-axis ranges from 10^{-2} to 10^{-1} . Two lines are plotted: a blue line for 'irányított' (directed) networks and an orange line for 'irányítatlan' (undirected) networks. Both lines show a negative correlation, with the undirected network having higher transitivity for a given number of peaks.

log(csúcsok száma)	log(transzitivitás) (irányított)	log(transzitivitás) (irányítatlan)
1	~0.03	~0.08
10	~0.015	~0.04
100	~0.007	~0.02

mondat	sugár, r	átmérő, d	center	periféria		
	$\min e_v$	$\max_v e_v$	$\{v \mid e_v = r\}$	$\{v \mid e_v = d\}$		
	\rightarrow	\leftarrow				
100	11	7	23	13	$\{., \text{!}, ?\}$	{nádcukorból}
1k	9		19		$\{., \}$	{Megadható, two}

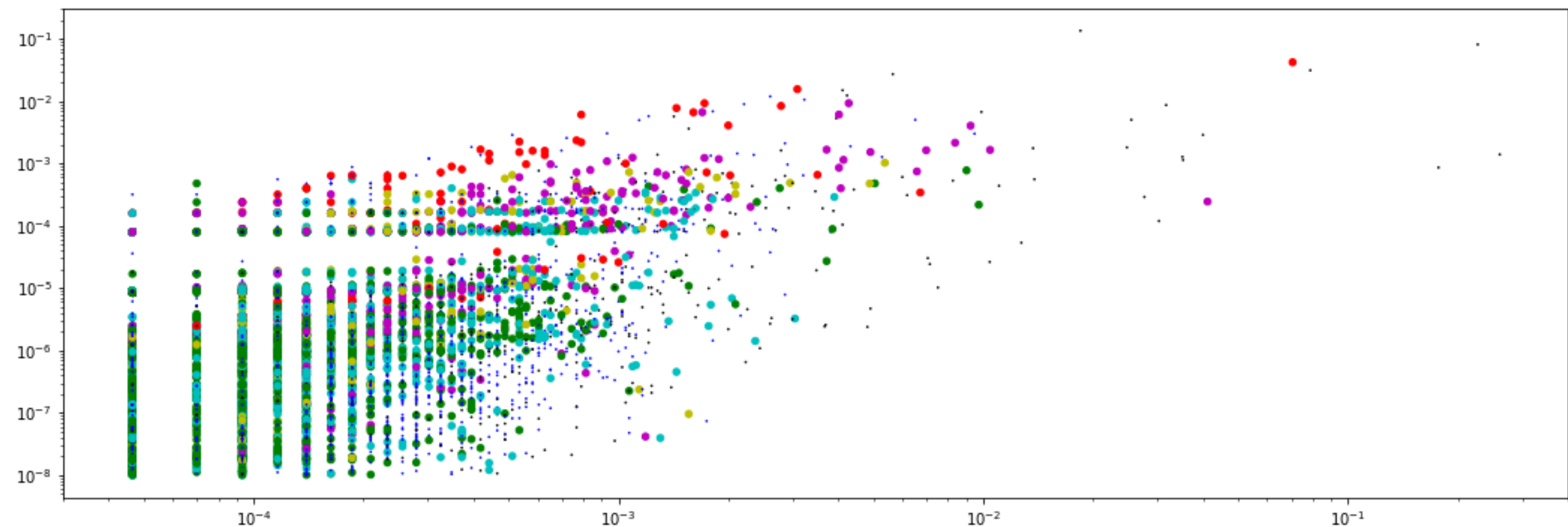
HITS, hyperlinkindukált témaeresés (*Hyperlink-Induced Topic Search*)

	pl.	miért fontos
hub	index.hu , vajdasag.lap.hu	linkek
tekintély (<i>authority</i>)	http://www.nytud.hu/oszt/korpusz/	tartalom

- kölcsönös definíció
- számítása iteratív
 - tetszőleges inicializáció
 - majd minden iterációban \rightarrow

$$\begin{aligned} h(v_1) &= \Sigma\{a(v_2) \mid \langle v_1, v_2 \rangle \in E\} \\ u(v_2) &= \Sigma\{a(v_1) \mid \langle v_1, v_2 \rangle \in E\} \\ u & /= \sum_v u(v)^2 \\ a & /= \sum_v a(v)^2 \end{aligned}$$

4. ábra A nagyobb, piros ponttal jelölt kötőszavak balra fent (magasabb authority), a nagyobb, zöld ponttal jelölt igék jobbra lent (alacsonyabb authority) helyezkednek el a fokszám (gyakoriság) vs authority grafikonon.



5. ábra Tekintély szófajok szerint: **kötőszók**, **igék**, **határozók**, **melléknevek**, és **számnevek**. 10 K mondat, csak a $> 10^{-8}$ tekintélyű szavakat ábrázoltuk.

- TextRank (Mihalcea and Tarau, 2004), kulcsszókinyerés
results [...] are worse than results obtained with undirected graphs, which suggests that [...] there is no natural “direction”
- szemantikus hálók (Steyvers and Tenenbaum, 2005)
- trigram (Ferrer i Cancho and Solé, 2001)
- a skálafüggetlen-hípe kritikája (Willinger et al., 2009)

A köztiség (betweenness)	B közelség (closeness)
C sajátvektor-	D fok (itt gyakoriság)
E harmonikus	F Katz

- klikkek szófajok szerint?
- az élsúlyok skálázása távolságként
- irányított gráfok hatékony implementációja
- szemantika

A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.

G. Fagiolo. Clustering in complex directed networks. *Physical Review E*, 76(2):026107, 2007.

R. Ferrer i Cancho and R. Solé. The small world of human language. *Proceedings of The Royal Society of London. Series B, Biological Sciences*, 268:2261–2266, 2001.

R. Mihalcea and P. Tarau. Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, 2004.

Cs. Oravecz, T. Váradi, and B. Sass. The Hungarian Gigaword Corpus. In *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC2014)*, Reykjavík, 2014.

M. Steyvers and J. B. Tenenbaum. The large-scale structure of semantic networks: Statistical analyses and a model of semantic growth. *Cognitive science*, 29(1):41–78, 2005.

W. Willinger, D. Alderson, and J. C. Doyle. Mathematics and the internet: A source of enormous confusion and great potential. *Notices of the American Mathematical Society*, 56(5):586–599, 2009.

G. K. Zipf. *The Psycho-Biology of Language; an Introduction to Dynamic Philology*. Houghton Mifflin, Boston, 1935.