

*

On Regularity of Daily Distribution of Queries in Search Engine

(Sang-Gue Park)**, (Chan-Kyu Lee)***, (Kyung-Hyun Yoon)****,
(Seong-Hee Kim)*****, (Jun-Ho Lee)*****

가
 , Pareto Zipf . 2
 Pareto , 가 1.33 1.34
 Pareto 가 .

ABSTRACTS

In this paper we analyzed regularity of daily patterns of distribution of Queries coming from internet search engine. And then, we proposed a Pareto distribution and Zipf law for identifying the query distribution and applied them to daily queries on the search engine during 2 week. We found that there is some evidence that Pareto and Zipf laws can be applied to evaluate the regularity of daily patterns of distribution of queries in search engine. Those results can be used to provide a better understanding of the social interests and trends using the query distribution patterns.

daily query patterns, pareto model, regularity, search engine, zipf's law.

*
 ** (spark@cau.ac.kr)
 *** (leeck@cau.ac.kr)
 **** (khyoon@cau.ac.kr)
 ***** (seonghee@cau.ac.kr)
 ***** (joonho@naver.com)

1.

가 , 가 가 가 가

가 가 . 가

가 , 2007 6 6

가 75.5% 가

3,443 2.0%가 가

가 , 88.7% 가

가

(search engine) 가

가 가

가 가

가 1), 가

가 가

Zipf

가 . Silverstein et al.(1999) 6 9

가 2,000 9

가 . Spink et al.(2001)

1) 가 가

1997

2001

2

2,500
가

가

. Jansen, Spink,
Pedersen(2005) 2002

2,603

가

Silverstein

et al. (1999)

가

가

(2002, 2005)

1

가

가

가

가 (pattern)

가 1
 가 가 ,
 가 .

2. Pareto Zipf

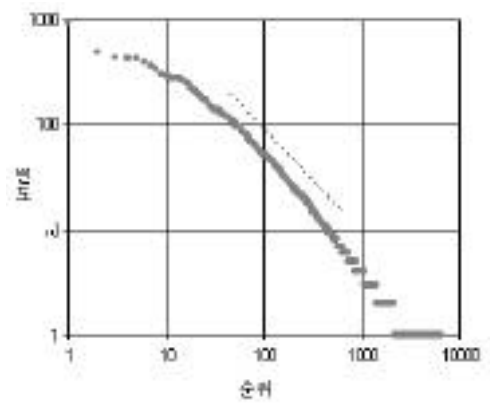
가
 가 (regularity)
 Pareto(1897)
 가 ,
 가

가 가 $Pr(X \leq x) = x^{-\alpha} (0 < \alpha < 1)$
 가 , 80-20

Zipf
 Pareto
 Pareto Hart &
 Oulton(1997) 가
 가 1.1~1.2
 (Goerge
 Kingsley Zipf, 1902-1950)
 가 2)

‘Zipf’ 가1 가
 =1 Zipf 가1
 Zipf 가1
 Zipf . 가1
 Pareto . Zipf
 가
 가
 1/3 ,
 가
 Zipf (Hart & Oulton, 1997),
 (Gabaix & Ioannides, 2003)
 Barasi & Albert(1999) WEB Page
 가
 (日間)
 Pareto

가1
 가1
 가1
 Zipf . Zipf
 Pareto
 가



< 1 > Zipf

$$F_i = \frac{\alpha}{R_i^\beta}, \quad i = 1, \dots, n, \quad \alpha, \beta > 0. \quad (1)$$

(1) R_i i-
 , F_i i-
 , Pareto . (1)

$$\ln(F_i) = \ln(\alpha) - \beta \ln(R_i). \quad (2)$$

(2) Pareto

3.

2003	1	1	14	2
	N			
			. N	
			N	

2)

가 , 가 Zipf
 . 2
 가 < 1 >
 가 1 1
 가 , 500
 Pareto 가 500
 < 1 >

			(, %)
2003.1.1()	4,081,960	203,550	110,657 (2.71)
2003.1.2()	5,471,997	237,259	139,081 (2.54)
2003.1.3()	5,406,911	236,736	136,248 (2.52)
2003.1.4()	4,638,271	221,052	119,654 (2.58)
2003.1.5()	4,075,251	211,048	105,887 (2.60)
2003.1.6()	5,551,299	243,639	138,215 (2.49)
2003.1.7()	5,805,736	243,895	278,919 (4.80)
2003.1.8()	5,724,836	244,701	184,606 (3.22)
2003.1.9()	5,567,749	244,313	146,964 (2.64)
2003.1.10()	5,541,086	242,444	138,071 (2.49)
2003.1.11()	4,833,728	223,557	121,095 (2.51)
2003.1.12()	4,198,654	215,682	108,231 (2.58)
2003.1.13()	5,987,183	250,013	142,990 (2.39)
2003.1.14()	5,836,792	250,135	138,941 (2.38)

가 . 가
 4.8% 가1 Zipf
 2.5% .
 가
 가 500 Pareto
 20 ~25 .
 500 가
 가
 가
 (2)
 < 2> Pareto
 Zipf
 < 2> 가 1 가
 가 80% 가
 < 2>

	ln ()		R ²
2003.1.1()	7.01538	1.34438	0.9738
2003.1.2()	7.20212	1.34854	0.9659
2003.1.3()	7.18124	1.34552	0.9665
2003.1.4()	7.09016	1.34235	0.9700
2003.1.5()	6.99815	1.33527	0.9727
2003.1.6()	7.19986	1.34344	0.9650
2003.1.7()	7.20158	1.34311	0.9657
2003.1.8()	7.20597	1.34322	0.9650
2003.1.9()	7.18029	1.33945	0.9655
2003.1.10()	7.19954	1.34468	0.9667
2003.1.11()	7.11344	1.34568	0.9707
2003.1.12()	7.01347	1.33487	0.9723
2003.1.13()	7.25226	1.34726	0.9636
2003.1.14()	7.21915	1.34181	0.9641

10

가 . , 가 가 가
 가 . , 가

(2006)

. < 2> -
 Zipf
 < 2> < 1> < 2>
 - Pareto

Zipf

Zipf
 (heavy tail)

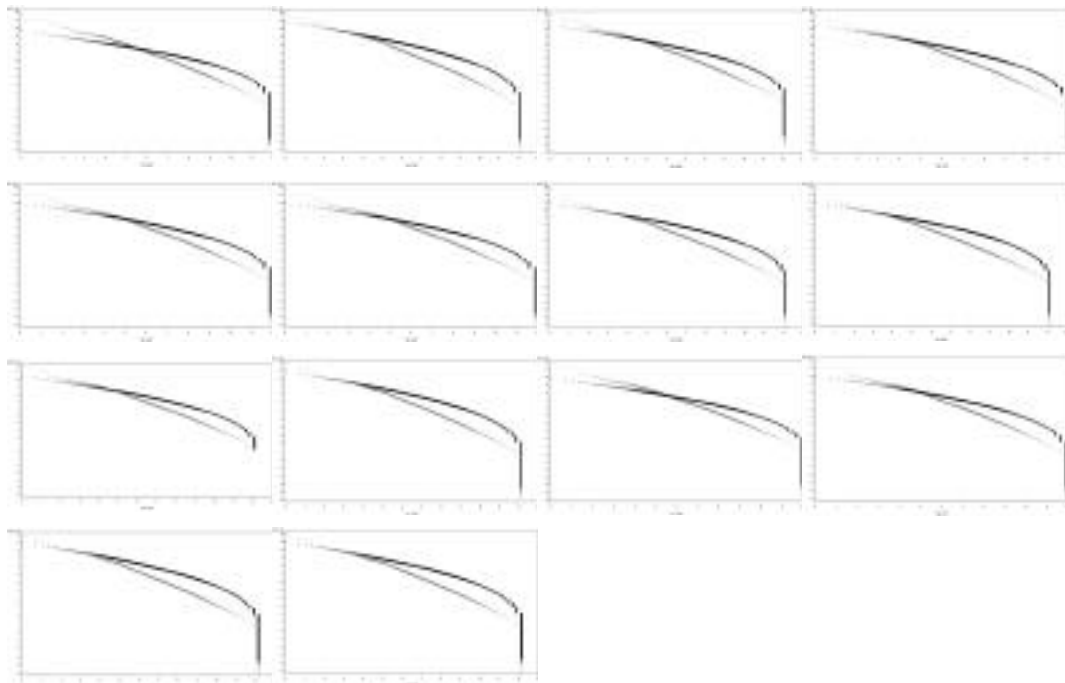
. 1 1 1 11

Zipf

Zipf

Zipf

Zipf



< 2> Zipf 2 Zipf

Zipf

Zipf
 (2) < 3>
 < 2> < 3> 4.
 5% 가
 , 가1 가 가
 가
 가 Pareto Zipf
 . 2
 가1 Zipf - Pareto
 .1 1 , 가1.33
 1 11 가1.1 1.34 Pareto

< 3>

	ln ()		R ²	
2003.1.1()	13.3666	1.1283	0.9978	9,001
2003.1.2()	12.4784	1.0340	0.9989	7,490
2003.1.3()	13.1278	1.0882	0.9980	12,384
2003.1.4()	12.8120	1.0709	0.9981	8,193
2003.1.5()	12.9540	1.0894	0.9978	8,740
2003.1.6()	12.1852	1.0068	0.9996	7,289
2003.1.7()	12.6255	1.0427	0.9991	7,194
2003.1.8()	12.3239	1.0175	0.9997	7,153
2003.1.9()	12.4915	1.0317	0.9989	7,237
2003.1.10()	12.2362	1.0115	0.9994	7,385
2003.1.11()	13.4920	1.1265	0.9976	16,270
2003.1.12()	12.9539	1.0870	0.9979	8,745
2003.1.13()	12.0710	0.9941	0.9999	6,952
2003.1.14()	12.2238	1.0069	0.9995	7,026

가 1 . 2 가

Pareto 1 Pareto Zipf 2 가

가 1 Zipf 가

- 가

가 Zipf 가

, Zipf Zipf 가

- Zipf 가 1.1

- Zipf (日) ,

. Pareto , 가

가 1 Zipf

, (日間) ,

Zipf

. 가

2003 , .

- , , 2005.
 ,
 39(2):146-160.
 , 2002.
 ,
 19(3):111-122.
 , . 2003.
 20(2)
 : 28-41.
 . 2006. Zipf
 ,
 22(2):275-299.
 , 2007. 2007
- Basasi, G. and R. Albert. 1999. "Emergence of scaling in random networks". *Science*, 286:509-512.
- Champernowne, D. 1953. "A model for income distribution". *Economic Journal*, 63:318-351.
- Gabaix, X. 1999. "Zipf's law of city sizes: An explanation". *Quarterly Journal of Economics*, 114(3): 739-767.
- Gabaix, X. 2000. "Zipf's law and the growth of cities". *American Economic Review*, 89(2): 129-132.
- Hart, P.E. and N. Oulton. 1997. "Zipf and the size distribution of firms". *Applied Economic Letters*, 4:205-206.
- Jansen, B. J., A. Spink, A., & J. Pedersen. 2005. "A temporal comparison of Alta Vista Web searching". *Journal of the American Society for Information Science and Technology*, 56(6): 559-570.
- Reed, W. 2001. "The Pareto, Zipf and other power laws". *Economic Letters*, 49:453-457.
- Silverstein, C., M. Henzinger, H. Marais and M. Moricz. 1999. "Analysis of a very large web search engine query log". *SIGIR. Forum* 33(1):6-12.
- Spink, A. et al. 2001. "searching the web: the public and their queries." *Journal of the American Society for information science and Technology*, 52(3):. 226-34.
- Zipf, G.K. 1949. *Human Behavior and the Principle of Least Effort*. Reading: Addison Wesley.

