

Meybodi@ce.aku.ac.ir

Tashakori@noavar.com

.....

... ([])

[]

[]

[]

[]

:

[]

[]

:

text retrieval
index
query
descriptive cataloging
subject cataloging

exhaustivity
 specificity
 recall
 precision

()

[]

[]

()" "

MS WORD" " PE2" " "

hard copy
convert

zipf

()

:[]

Frequency.rank ≈ Constant

()

zipf (... " " " " " ") []

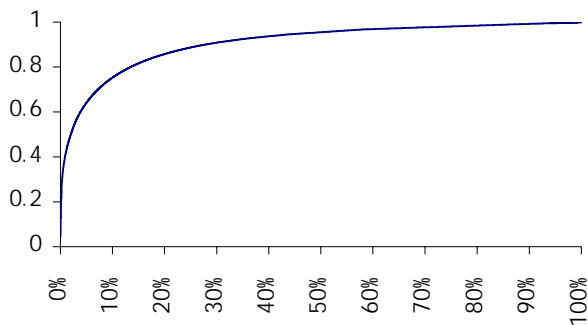
" " zipf []

" "

Zipf

(=)

(R)	(F)	R × F/N	(R)	(F)	R × F/N
21	337	0.100487	1	3007	0.042697
22	331	0.103398	2	2475	0.070286
23	311	0.101566	3	2030	0.086473
24	309	0.105301	4	1640	0.093146
25	295	0.104718	5	1619	0.114942
26	288	0.106323	6	1338	0.11399
27	286	0.109645	7	1077	0.107047
28	278	0.110526	8	1053	0.119613
29	272	0.112002	9	1020	0.130348
30	227	0.096696	10	960	0.136311
31	218	0.095958	11	846	0.132137
32	218	0.099053	12	792	0.134948
33	218	0.102148	13	660	0.121828
34	208	0.100416	14	528	0.10496
35	206	0.102376	15	525	0.111818
36	204	0.104278	16	505	0.114729
37	201	0.105599	17	448	0.10814
38	195	0.105215	18	420	0.107345
39	180	0.099678	19	394	0.106294
40	177	0.10053	20	391	0.111037



()

stop list, negative dictionary

(negative dictionary stop list)

'is' 'the' 'a' .[]

%

% " "

.[]

)

(

(" ") (" ")

" "

" "

" "

(" ")

" " " "

"

"

.[]

B A

B A

:[]

B A

(**A ∪ B**)

"B A"

•

B A

(**A ∩ B**)

"B A"

•

B A

(**A - B**)

"B A"

•

Boolean

" " ()

()

" " " " " "

()

[]

[]

n-gram

n-gram

[]

C A Ax C → Ay C

³ successor variety

¹ table lookup
² affix remove

⁴ recoding

.....

[]

.....

.....

.....

:

[]

[]

..

..

..

[]

..

..

--

..

[]

[]

(" " " " " " " ")

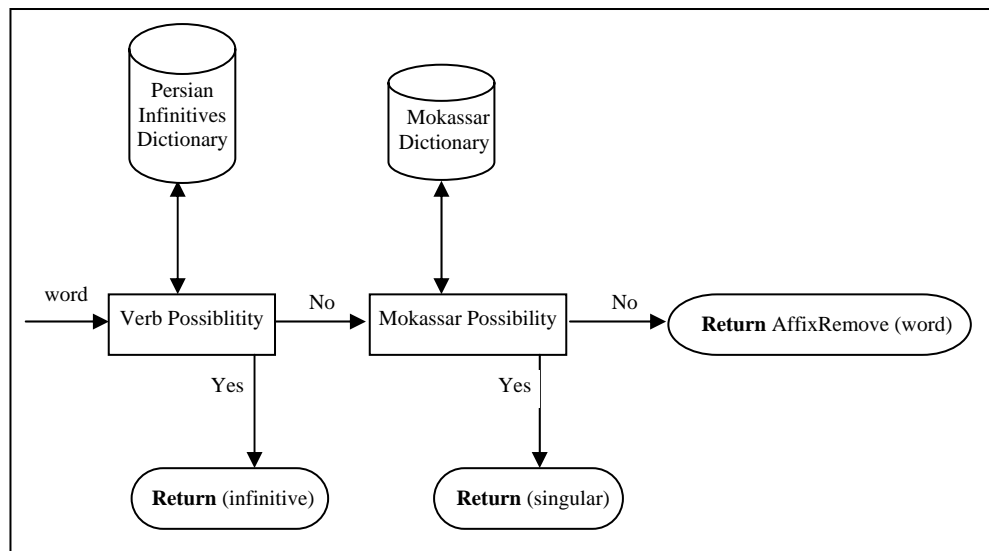
..

() ()
 () ()
 " " ()

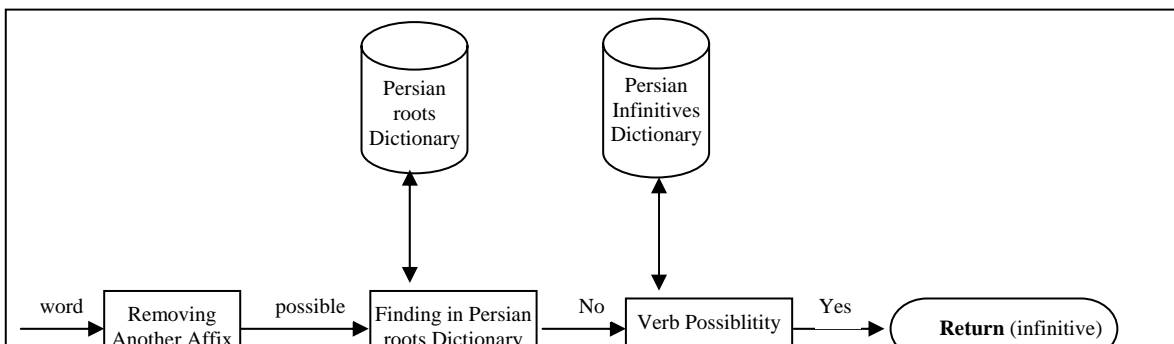
Stem(word) word

" "

Stem(word):



AffixRemove(word):



" "

[]

[]

[]

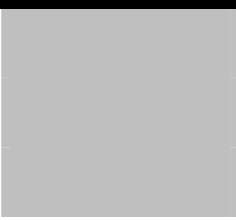
:

" "

()

¹ understemming
² overstemming
³ light stemmer
⁴ heavy stemmer
⁵ correctness

" "



% / % /

" "

64Mb PIII 500 Mhz

/

-

" "

" "

" "

" "

" "

" "

	" "	" "
	0.3595258	0.8974702
	0.5421372	0.8397220

" "

" " "

" "

" "

%

" "

%

" " " "

- [2] Chen H., Schatz B., Ng T., Martinez J., Kirchoff A., and Lin C., "A Parallel Computing Approach to Creating Engineering Concept Spaces for Semantic Retrieval: The Illinois Digital Library Initiative Project", *IEEE Trans. on Pattern Analysis and Mach. Intell.*, Vol. 18, No. 8, pp. 771-782, 1996.
- [3] Hmeidi I., Kanaan G., and Evens M., "Design and Implementation of Automatic Indexing for Information Retrieval with Arabic Documents", *Journal of the American Society for Information Science*, Vol. 48, No. 10, pp. 867-881, 1997.
- [4] Salton G. and Mc Gill M. J., *Introduction to Modern Information Retrieval*, Mc Graw Hill, New York, 1983.
- [6] Rijsbergen C.J. Van, *Information Retrieval*, <http://www.dcs.gla.ac.uk/Keith>
- [7] Luhn H. P., "A Statistical Approach to Mechanized Encoding and Searching of Literary Information", *IBM Journal of Research and Development*, Vol. 1, No. 4, pp. 309-317, 1957.
- [8] Salton G., *Automatic Text Processing*. Reading, Mass.: Addison-Wesley, 1989.

¹ general accuracy
² length truncation

- [10] Paice C. D., "An Evaluation Method for Stemming Algorithms", *Proc. of the seventeenth annual international ACM-SIGIR conference on Research and development in information retrieval*, pp. 42-50, 1994.
- [11] PORTER M. F., "An Algorithm for Suffix Stripping", *Program*, vol. 14, no. 3, pp. 130-137, 1980.
- [12] PAICE, C. D., "Another Stemmer", *ACM-SIGIR Forum*, vol. 24, no. 3, pp. 56-61, 1990.
- [13] DAWSON J., "Suffix Removal and Word Conflation", *ALLC Bulletin*, Michemas, pp. 33-46, 1974.
- [14] LOVINS J. B., "Development of a Stemming Algorithm", *Mechanical Translation and Computational Linguistics*, vol. 11, no. 1-2, pp. 22-31, 1968.
- [15] Frakes W. B., *Stemming Algorithms*, <http://matrix.nbu.bg/books/books/book5/chap08.htm>

[]

[]

[]

[]

[]

[]

[]

- [22] Lennon M., Pierce D. S., Tarry B. D., and Willett P., "an Evaluation of Some Conflation Algorithms for Information Retrieval", *Journal of Information Science*, no. 3, pp. 177-183, 1981.

	()		
	()		
	()		

	()		
	() ()		