

()

موضوع بندی متن های مترجم

**

*

*

**

(

)

(

)

:

۱. Dense Text

Topic Specific of the Dense Documents

N. Ghasem Aghaee and G. Sarrafan

Department of Computer Engineering, University of Isfahan

Abstract

This paper investigates text documents regarding their topic density. It has divided them into two groups: dense and sparse documents. Dense documents are texts with a wide domain of topics. They have a high topic density (for example religious books, encyclopedia, magazine archives, etc). We have shown that a) traditional methods can not be used for topic

1. Triples

()

/...

specific of dense texts, and b) we can benefit from employing the efficiency of the proposed method (Nasir) for dense texts.

In this research, we have used dependency relations, paths, triple databases and statistical text processing methods to extract important words and to insert them into a clustering index. Also a method was described to find the reference of pronouns in dense texts.

In addition, based on the suggested methods, a prototype system called Nasir was implemented. The result of the implementation on Persian dense texts shows that the quality of indexing and searching improved significantly.

Keywords: Text processing, Dense text, Topic specific, Dependency relations, Pronoun referencing

« »
« »
:
:
()
()
:
()

TFIDF

1. Term Frequency Inverse Document Frequency

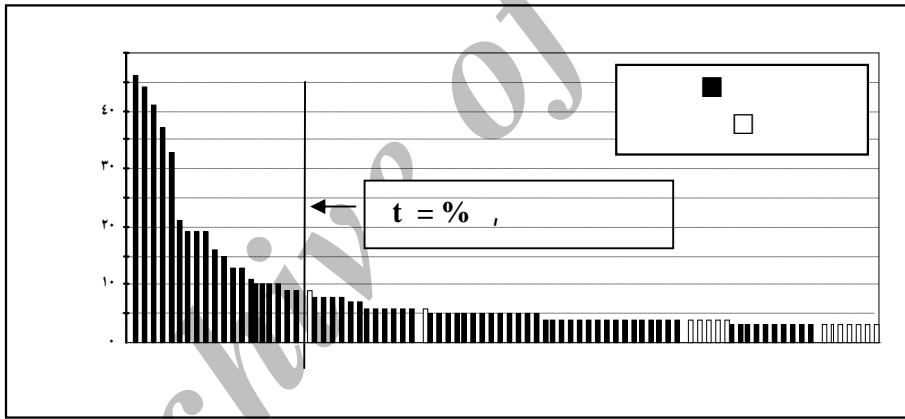
/...

: TFIDF

()

TFIDF

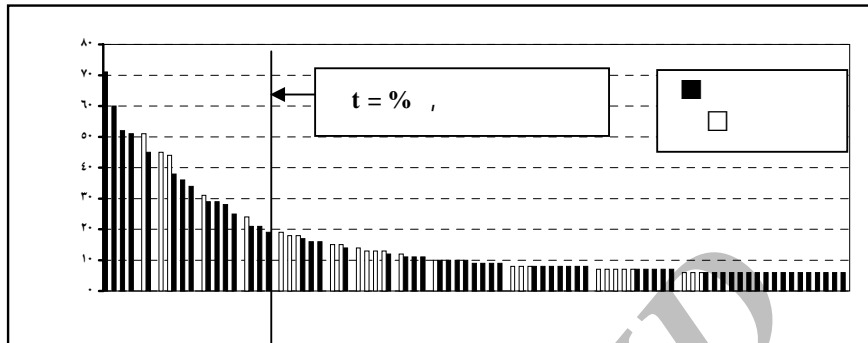
).
()



:

Archive of SID

/



(t)

/	/	()
/	/	
/		
(/)	()	()

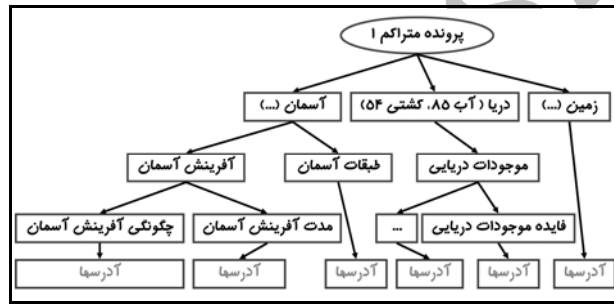
(.)

/...

()

()

()



:

()

generality(t) t

. Thesaurus

$$\text{generality}(t) = \frac{\sum_d \frac{f_d(t)}{L_d}}{N}$$

$$f_d(t) = \sum_{t=1}^N \frac{f_d(t)}{L_d}$$

$$\text{generality}(t) = \frac{\sum_d \frac{f_d(t) \cdot \text{density}(d)}{L_d}}{N}$$

$$\text{density}(d) = \frac{N_{\text{imp_term}}(d)}{L_d}$$

$$N_{\text{imp_term}}(d) = \sum_i^{L_d} f_{t_i} | w(t_i) > w_{\text{threshold}}$$

$$\text{generality}(t) = \frac{\sum_d \frac{f_d(t) \cdot \text{density}(d)}{L_d}}{N}$$

1. Iterative

/...

$$w(t) = \frac{1}{generality(t)}$$

w_threshold

w(t)

w_threshold

generality

d

generality

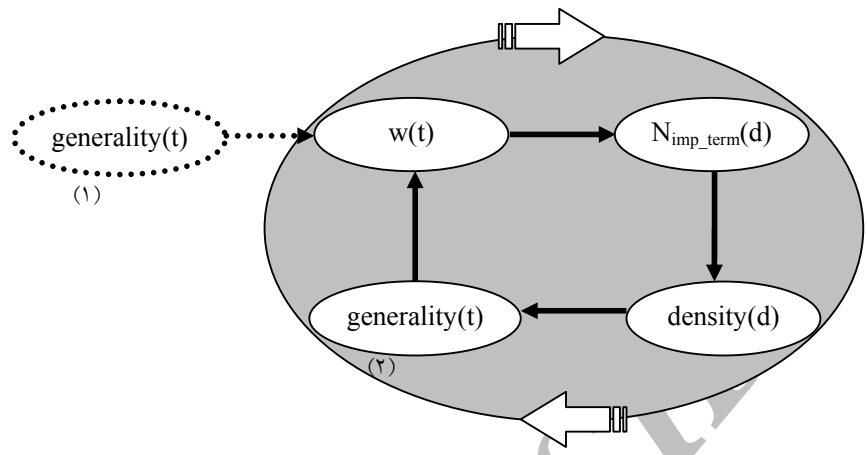
generality

generality

Archive of SID

w(t)

generality(t)



:

t

$w(t)$

Δw

$w(t)$

t

$+ \Delta w$

:()

$$\Delta w = \frac{w_threshold}{\gamma \times (\log_{\gamma}(N_{keyword} + 1))}$$

$N_{keyword}$

Δw

()

/...

()

$ti() \quad set = [t_1, t_2, \dots, t_n]$

$\forall ti \in Set \mid imp_term(ti)$

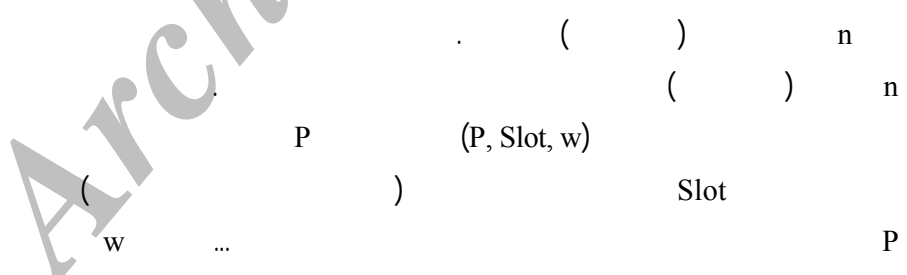
→

$child(set, ti)$

$ti \quad (« \dots »)$
 $(« \dots » « \dots »)$

$(« \dots ») Set$
 ti

Set



- v. Stack
- r. Filters

/

Slot

(. ()).

.() (Y X)

: () ()

$$mi(p, Slot, w) = \log \left(\frac{|p, Slot, w| \times |*, Slot, *|}{|p, Slot, *| \times |*, Slot, w|} \right)$$

: p « »

« . »
N: : V ← → V: : N : P
. Y X :

(P , SlotX , « »)

X P

. « »

1. Slot

/...

« »

« »:

« »

« »

()

(/)

:

(Find)

« » « »

- ۱. Precision
- ۲. Recall

/

:

« »

« »

(« »)

: « »

« »

Archive of SID

/...

()

:

			()

Archive of SID

)

(«

»

«

»

.()

« »

().

() TFIDF

()

Minipar

Archive of SID

١. Parser
٢. Minipar: www.cs.ualberta.ca/~lindek/minipar.htm.

/...

- . D. Lin and P. Pantel, "DIRT: Discovery of Inference Rules from Text,"
In Proceedings of the ACM SIGKDD Conference on Knowledge
Discovery and Data Mining,
- . J. J. Lin, "Indexing and Retrieving Natural Language Using Ternary
Expressions", Master's Thesis, Massachusetts Institute of Technology,
- . Kjersti Aas and Line Eikvil, "Text Categorization: a Survey", www.citeseer.nj.nec.com, June

Archive of SID