

:	:	2022-07-26
:	:	
:	:	0%
:	:	0.00
:	:	0.00

1 - 2022-07-26 15:27 -

- clipboard-202207261521-xhtfb.png (가) 가 .
- clipboard-202207261522-jsmd1.png (가) 가 .
- clipboard-202207261526-fszqw.png (가) 가 .

TF-IDF (Term Frequency-Inverse Document Frequency)

(1)

- 가 가
- , " " 가
- , ,

(2) TF (Term Frequency)

$$tf(t,d) = 0.5 + 0.5 \cdot \frac{freq(t,d)}{\text{(문서 내 단어들의 } freq(t,d) \text{ 값 중 최대 값)}}$$

- TF 가
- Term , 가
- 가 가

(3) IDF (Inverse Document Frequency)

$$idf(t,D) = \log\left(\frac{|D|}{1 + |\{d \in D : t \in d\}|}\right) = \log\left(\frac{\text{전체문서의 수}}{1 + \text{단어 } t \text{가 포함된 문서의 수}}\right)$$

- IDF , DF
- , DF
- Term 가
- (N)가 log , 가
- 가 , 가 0 1

(4) Score

$$tfidf(t,d,D) = tf(t,d) \times idf(t,D)$$

- TF-IDF Score TF IDF

1) TF

1	0	0	0	1
2	0	0	0	0
3	1	1	2	0
4	1	0	0	0

- (DTM : Document-Term Matrix)

2) IDF

	IDF
	$\ln(4/(2+1))=0.288$
	$\ln(4/(1+1))=0.693$
	$\ln(4/(2+1))=0.288$
	$\ln(4/(1+1))=0.693$

- (N)가 4 ln 4

- 가 2 2

3) TF*IDF

1	0	0	0	0.693
2	0	0	0	0
3	0.288	0.693	0.576	0
4	0.288	0	0	0

- 2) IDF 가

- IDF 0.288 3 2 0.288*2

#2 - 2022-07-27 19:04 -

- () #6520 #6552() .

clipboard-202207261521-xhtfb.png	15.9 KB	2022-07-26
clipboard-202207261522-jsmd1.png	23.7 KB	2022-07-26
clipboard-202207261526-fszqw.png	8.44 KB	2022-07-26