

# TF - IDF

Q1) Doc - Database.

\* Corpus D :

	man
D1 : Information Retrieval System	3
D2 : Information Storage	2
D3 : Digital Speech Synthesis System	4
D4 : Speech Filtering	2
D5 : Speech retrieval	2

(i) TF ~~TF~~ value table: (Term Freq. Matrix)

		Document / Term.				
Words.		D1	D2	D3	D4	D5
w1	Information	1/3	1/2	0	0	0
w2	Retrieval	1/3	0	0	0	1/2
w3	system	1/3	0	1/4	0	0
w4	Storage	0	1/2	0	0	0
w5	Digital	0	0	1/4	0	0
w6	Speech	0	0	1/4	1/2	1/2
w7	Synthesis	0	0	1/4	0	0
w8	Filtering	0	0	0	1/2	0

\* Formula :

$$TF(w, d) = \frac{\text{Occurrence of } w \text{ in } d}{\text{total no. of } w \text{ in } d}$$

-D]

IDF formula:  $IDF(w_1 \text{ in } D) = \log\left(\frac{N}{df_1}\right) = \log\left(\frac{5}{2}\right)$

↳ Corpus D

② TF-IDF (Term Doc. Freq)

Formula:  $w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$

$tf_{i,j}$  = no. of occurrences of  $i$  in  $j$

$df_i$  = no. of documents containing  $i$

$N$  = total no. of documents

TF-IDF( $w_1$  in  $D_1$ ) =  $0.33 \times \log\left(\frac{5}{2}\right) = 0.1313$

TF-IDF( $w_2$  in  $D_1$ ) =  $0.33 \times \log\left(\frac{5}{2}\right) = 0.1313$

TF-IDF( $w_3$  in  $D_1$ ) =  $0.33 \times \log\left(\frac{5}{2}\right) = 0.1313$

⋮

② Find IDF for words:

Formula:  $IDF(w_i, D) = \log\left(\frac{N}{df_i}\right)$

$IDF(w_1, D) = \log\left(\frac{5}{2}\right) = 0.39$

$IDF(w_2, D) = \log\left(\frac{5}{2}\right) = 0.39$

$IDF(w_3, D) = \log\left(\frac{5}{2}\right) = 0.39$

$IDF(w_4, D) = \log\left(\frac{5}{1}\right) = 0.69$

$IDF(w_5, D) = \log\left(\frac{5}{1}\right) = 0.69$

$IDF(w_6, D) = \log\left(\frac{5}{3}\right) = 0.22$

$IDF(w_7, D) = \log\left(\frac{5}{1}\right) = 0.69$

$IDF(w_8, D) = \log\left(\frac{5}{1}\right) = 0.69$

Teacher's Signature:.....

③ Find TF-IDF for word = "Speech"

TF-IDF ("Speech" in D) =

Table TF-IDF for word = Speech.

Document	TF(speech)	IDF(speech)	TF-IDF
D1	0	0.22	0
D2	0	0.22	0
D3	1/4	0.22	0.055
D4	1/2	0.22	0.11
D5	1/2	0.22	0.11

TF-IDF formula:

$$\text{TF-IDF}(\text{word in Doc}) = \text{tf} * \log \left( \frac{N}{\text{df}_i} \right)$$

$$\begin{aligned} D1 &\rightarrow 0 * \log 0.22 = 0 \\ D2 &\rightarrow 0 * \log 0.22 = 0 \\ D3 &\rightarrow 0.25 * \log 0.22 = 0.055 \\ D4 &\rightarrow 0.5 * \log 0.22 = 0.11 \\ D5 &\rightarrow 0.5 * \log 0.22 = 0.11 \end{aligned}$$

So the TF-IDF values for the term "speech" in 2 docs are approx 0.11 for D4, D5 & for D3 it's 0.055. This means that speech is most relevant to these 3 docs. Among 3 two are more relevant to the term 'Speech' because they have the highest TF-IDF score.



Q 2) Explain.

A. Synonym:

A synonym is a word or phrase that has a similar or identical meaning to ~~other~~ another word or phrase in the same language.

-eg: happy & joyful

B. Antonym:

An Antonym is a word that has the opposite meaning of another word.

-eg: hot & cold

C. Hyponym:

It is a word whose meaning is a subset of a more general word.

-eg: Rose & flower

D. Hyperonym:

It is a word that represents a more general category or superclass of words.

-eg: Fruit & apple

E. Meronym:

It is a word that represents a part or a component of a larger whole.

-eg: wheel & car.

a) Couch - sofa

Synonym

Couch & sofa are synonyms because they both refer to a similar piece of furniture designed for seating.

b) Car - wheel.

Meronymy:

Wheel is a part of a car, so it's a meronymic relationship.

c) Meal - Breakfast:

Hyponym

Breakfast is a specific type of meal, so it's a hyponymic relationship.

Breakfast is a subset of meal.

d) I left my heart - and my suitcase.

No direct lexical relationship.

No direct lexical relation b/w heart & suitcase.

they do not represent typical word relationship

e) Mammal - Dog:

Hyponym

Dog is a specific type of mammal, so it's a hyponymic relation.

Dogs belongs to category mammals.

Q3)

- A: Jupiter is the largest planet  
 B: Mars is the fourth planet from the earth

\* Vector Form

a)

① Traditional word count method.

Words	A	B
Jupiter	1	0
is	1	1
the	1	2
largest	1	0
planet	1	1
Mars	0	1
Fourth	0	1
from	0	1
earth	0	1

$$\therefore \text{Vector A} = [1, 1, 1, 1, 1, 0, 0, 0, 0]$$

$$\text{Vector B} = [0, 1, 2, 0, 1, 1, 1, 1, 1]$$

② estimate similarity b/w doc A & B using <sup>similarity</sup> cosine Formula:

$$\text{Similarity (A, B)} = \frac{\sum A_i \times B_i}{\sqrt{\sum A_i^2} \times \sqrt{\sum B_i^2}}$$

$$\therefore A \cdot B = (1 \times 0 + 1 \times 1 + 1 \times 2 + 1 \times 0 + 1 \times 1 + 0 \times 1 + 0 \times 1 + 0 \times 1 + 0 \times 1)$$

$$= 0 + 1 + 2 + 0 + 1 + 0 + 0 + 0 + 0$$

$$= 4$$

$$|A| = \sqrt{5}$$

$$|B| = \sqrt{10}$$

Teacher's Signature: .....



$$\text{IDF} = \log(2)$$

$$\therefore \text{Similarity}(A, B) = \frac{4}{\sqrt{5} \times \sqrt{10}} = \frac{4}{\sqrt{50}} = \frac{4}{5\sqrt{2}} = 0.56$$

b) TF-IDF method. cosine similarity.

words	TF(A)	TF(B)	IDF	TFIDF(A)	TFIDF(B)
Jupiters	1/5	0	$\log(2/1) = 0.30$	0.06	0
Is	1/5	1/8	$\log(2/2) = 0$	0	0
the	1/5	2/8	$\log(2/2) = 0$	0	0
largest	1/5	0	$\log(2/1) = 0.30$	0.06	0
planet	1/5	1/8	$\log(2/2) = 0$	0	0
mass	0	1/8	$\log(2/1) = 0.30$	0.06	0.0375
fourth	0	1/8	$\log(2/1) = 0.30$	0.06	0.0375
from	0	1/8	$\log(2/1) = 0.30$	0.06	0.0375
earth	0	1/8	$\log(2/1) = 0.30$	0.06	0.0375

$$\text{Vector A} = [0.06, 0, 0, 0.06, 0, 0, 0, 0, 0]$$

$$\text{Vector B} = [0, 0, 0, 0, 0, 0.0375, 0.0375, 0.0375, 0.0375]$$

$$\text{Similarity}(A, B) = \frac{A \cdot B}{|A| \times |B|}$$

$$\therefore A \cdot B = 0 \quad |A| = 2 \quad |B| = 4$$

$$\text{Similarity}(A, B) = \frac{0}{2 \times 4} = \frac{0}{8} = 0$$

Teacher's Signature:.....

c) Traditional word count method yields a higher cosine similarity of approx. <sup>0.56</sup> ~~0.6761~~ indicating that the 2 docs share more common words. However it doesn't consider the importance of words.

TF-IDF method takes into account the importance of words. It gives a lower cosine similarity of approx 0, indicating that the 2 docs are less similar when considering the importance of words.

d) Other similarity measures.

1. Jaccard Similarity.

measures the similarity b/w 2 sets by comparing the size of their intersection & size of their union.

$$\text{Jaccard Similarity}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

2. Euclidean Distance:

It measures the straight-line distance b/w 2 points in a multidimensional space.

$$\text{Euclidean Distance}(A, B) = \sqrt{\sum_{i=1}^n (A_i - B_i)^2}$$

3. Pearson Correlation Coefficient:

measures the linear relationship b/w 2 vectors.

$$\text{Pearson Corr Coef}(A, B) = \frac{\sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B})}{\sqrt{\sum_{i=1}^n (A_i - \bar{A})^2} \sqrt{\sum_{i=1}^n (B_i - \bar{B})^2}}$$

$$\sqrt{\sum_{i=1}^n (A_i - \bar{A})^2} \sqrt{\sum_{i=1}^n (B_i - \bar{B})^2}$$

Teacher's Signature: .....