

## TFIDF :

It is combination of two things.

- ① Term frequency and ② Inverse documents frequency

① Term frequency =

freq. of term 't' in documents

Total no. of words in that documents.

"Data Science is a combination of ml and dl"

data science combination machine learning  
deep learning.

Teacher's Signature:.....

term freq. of "learning" =  $\frac{2}{1+1} = 1$

$\Rightarrow \frac{2}{7} = 0.2842$

T.f (Data) =  $\frac{1}{7} = 0.1428$

Document frequency  $\Rightarrow$

$\frac{\text{No. of documents containing 't'}}{\text{total no. of documents}}$

① NLP needs Deep learning and machine learning

② Data science contains lots of things like NLP

document freq of (NLP) =  $\frac{2}{2} = 1$

D.f of "learning" =  $\frac{1}{2} = 0.5$

Inverse document frequency:

IDF  $\Rightarrow \log \left( \frac{1}{\text{Document Frequency}} \right)$

$\Rightarrow \log \left( \frac{\text{Total no. of Document}}{\text{no. of Doc. containing 't'}} \right)$



$$\text{IDF} = \log\left(\frac{1}{1}\right) = \log(1) = 0$$

(NLP)

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

$$\text{TFIDF} \Rightarrow \text{TF} \times \text{IDF}$$

$$\text{TFIDF}(\text{learning}) = \frac{2}{1} \times \log(2)$$

exa!

- ① learning nlp.
- ② data science combinations deep learning machine learning
- ③ nlp needs machine learning deep learning.

Unique words:

in this case every unique word treated as a feature.

learning	nlp	data	science	combinations	machine	deep	needs
1)							
2)							
3)							

No. of row will be 3.

Teacher's Signature:.....



1st Documents

(TF)

(IDF)

$$\text{TFIDF (learning)} = \frac{1}{2} \log\left(\frac{2}{1}\right)$$

$$\Rightarrow \frac{1}{2} \times 0.6931 = 0.3466$$

$$\Rightarrow \underline{\underline{0}}$$

$$\text{TFIDF (nlp)} \rightarrow \frac{1}{2} \log\left(\frac{3}{2}\right)$$

Doc 2

$$\text{TFIDF (data)} = \frac{1}{2} \log\left(\frac{3}{1}\right)$$

$$\text{TFIDF (science)} = \frac{1}{2} \log\left(\frac{3}{1}\right)$$

$$\text{TFIDF (combination)} = \frac{1}{2} \log\left(\frac{3}{1}\right)$$

$$\text{TFIDF (maclure)} = \left(\frac{1}{2}\right) \log\left(\frac{3}{2}\right)$$

$$\text{TFIDF (deep)} = \frac{1}{2} \log\left(\frac{3}{2}\right)$$

$$\text{TFIDF (needs)} = \frac{1}{6} \log\left(\frac{3}{1}\right)$$

$$\text{TFIDF (learning)} = \frac{2}{7} \log\left(\frac{3}{3}\right)$$

Doc 3

$$\text{TFIDF (nlp)} = \frac{2}{6} \log\left(\frac{3}{2}\right)$$

$$\text{TFIDF (needs)} = \frac{1}{6} \log\left(\frac{3}{1}\right)$$

Teacher's Signature



$$\text{TFIDF}(\text{machine}) = \frac{1}{6} \log\left(\frac{3}{2}\right)$$

$$\text{TFIDF}(\text{deep}) = \frac{1}{6} \log\left(\frac{3}{2}\right)$$

$$\text{TFIDF}(\text{learning}) = \frac{2}{6} \log\left(\frac{3}{3}\right)$$

In this case we are just checking the weightage of the word in their respective documents.

freq of word  $\uparrow \Rightarrow$  weightage of word  $\downarrow$

### Drawbacks of TFIDF

- ① Curse of dimension
- ② Index is not maintained.
- ③ It is not understanding actual meaning of the word.

exa!  
movie review

- ① Bad movie
- ② Gwd movie
- ③ fabulous movie
- ④ Awesome movie
- ⑤ Average movie

In this review movie is not important where as the word along with



## Similarities bet<sup>n</sup> Count Vectorizer and TFIDF

- ① Drawbacks of both are same.
- ② When we initiate their model, the parameters of the both are same.

### Diffs

Count Vectorizer  $\Rightarrow$  Count the word  
 TFIDF  $\Rightarrow$  check the weightage of word.

While model training the diff bet<sup>n</sup> the accuracy of both is around 1-2%, but in reality it is very high, so from where we get the good evaluation matrix value we will continue with it.

Sparse matrix: It is a output matrix of count vectorizer.

In this matrix most of the values are zero.

So we convert this matrix into array.

### Data leaking $\rightarrow$

In this case we 1st split the data after that we perform the operation of preprocessing on both training side and testing side separately.