

# TFIDF

TERM FREQUENCY - INVERSE DOCUMENT  
FREQUENCY

-Devendra Kumar Arya

# TFIDF

- is a numerical statistic that reflects how important a word is to a document in a collection or corpus(a large or structured set of texts)

- TFIDF is the product of *term frequency* and *inverse document frequency*. There are various ways to evaluate these two.

- 83% of text based recommender systems use tf-idf

- Term Frequency: No. of times a term occurs in a document

According to **Hans Peter Luhn(1957)**, the weight of a term in a document is proportional to its term frequency

# Term Frequency: a simple raw-count approach

-Term Frequency: a raw-count approach

the no of times a term  $t$  occurs in a document  $d$ :  $\text{tf}(t,d) = f(t,d)$

“The clouds are in the sky.”

Other Approaches:

1. Boolean:  $\text{tf}(t,d) = 1$  if  $t$  exists in  $d$ , 0 otherwise
2. Normalized term frequency:  $f(t,d) / (\text{No of words in } d)$
3. Logarithmic Scaled:  $\log(1 + f(t,d))$
4. Augmented Frequency: to prevent a bias towards longer documents

$$\text{tf}(t,d) = 0.5 + 0.5 \cdot \frac{f_{t,d}}{\max\{f_{t',d} : t' \in d\}}$$

# Inverse Document Frequency

-There may be words which occur very frequently in all the documents but are not that important like `the`

-*Inverse Document Frequency* is used to diminish weight of such terms and increase weight of terms that occur rarely

-Inverse Document Frequency: a measure of how much information the word provides

-**Karen Sparck Jones(1972)** calls this *term specificity*. Specificity is inversely proportional to the no of documents in which that term is present

$$\text{idf}(t,D) = \log (N/(1+\{\#d : d \in D \text{ and } t \in D\}))$$

N - total no of documents

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \cdot \text{idf}(t, D)$$

Note-Denominator  $\rightarrow 0$  as the same term appears in more documents. Reduces importance of commonly occurring words

# Meaning

-TF: The more frequent the term, the higher the score

-IDF: The more common the term, the lower the score

“The clouds are in the sky.”

-Since "the" is very common in every article, the score of IDF would be very low compared to “clouds” and “sky” which are more important terms for a document than “the”.

# Demo

# Conclusion

- Tfidf is a really popular technique for weighting the importance of the terms inside a collection of documents
- is used in Information Retrieval to rank results. One can also create a page ranking algorithm using this simple concept of TFIDF
- is used for extracting keywords on web pages

Sources:

<https://en.wikipedia.org/wiki/Tf-idf>

<https://nlpforhackers.io/tf-idf/>

<https://hackernoon.com/the-fastest-way-to-identify-keywords-in-news-articles-tfidf-with-wikipedia-python-version-baf874d7eb16>



Suggestions for improvements..