# TF - IDF ( Team frequency - Inverse Document frequency )

TO evaluate the importance of word in document or corpus
↓
set

① TF - IDF is a statistical Measure used to evaluate the importance of word in document relative to a collection of document

② It helps to weigh down common words that occur frequently across documents but are less informative, while giving more weights to rarel words that might be more significant for specific document

$$TF = \frac{No \ of \ rep \ of \ words \ in \ sentence}{No \ of \ words \ in \ sentence}$$

Measures how frequently a word appears in a document

$$IDF = \log_e \left( \frac{No \ of \ sentences}{Nof \ of \ sentences \ containing \ the \ word} \right)$$

Measures how important a word in given corpus

IDF ↑ - zero

$$\boxed{Tf - IDF = Term \ frequue \times Inverse \ Document \ frequeny}$$

Example

Doc1 : " Cat eats fish "

Doc2 : " dog eats fish "

Doc3 : " dog chases cat "

Term frequency

| Term | DOC-1 (cat eats fish) | Doc 2 (dog eats fish) | Doc3 (dog chases cat) |
|---|---|---|---|
| Cat | $\frac{1}{3}$ | 0 | $\frac{1}{3}$ |
| eats | $\frac{1}{3}$ | $\frac{1}{3}$ | 0 |
| fish | $\frac{1}{3}$ | $\frac{1}{3}$ | 0 |
| dog | 0 | $\frac{1}{3}$ | $\frac{1}{3}$ |
| chases | 0 | 0 | $\frac{1}{3}$ |

**Step 2** — Calculate inverse Document Frequency
[#g document containing word]

| | | |
|---|---|---|
| Cat | 2 | $\log\left(\frac{3}{2}\right) = 0.18$ |
| eats | 2 | $\log\left(\frac{3}{2}\right) = 0.18$ |
| fish | 2 | $\log\left(\frac{3}{2}\right) = 0.18$ |
| dog | 2 | $\log\left(\frac{3}{2}\right) = 0.18$ |
| chases | 1 | $\log\left(\frac{3}{1}\right) = \boxed{0.40}$ |

**TF - IDF**

**DOC1**

$$TF-IDF(cat) = \frac{1}{3} \times 0.18 = 0.06$$

$$eats = \frac{1}{3} \times 0.18 = 0.06$$

$$fish = \frac{1}{3} \times 0.18 = 0.06$$

**DOC 2**

$$dog = \frac{1}{3} \times 0.18 = 0.06$$

$$cats = \frac{1}{3} \times 0.18 = 0.06$$

$$fish = \frac{1}{3} \times 0.18 = 0.06$$

**DOC3**

$$\frac{TF \times IDF}{}$$

$$dog = \frac{1}{2} \times 0.18 = 0.06$$

$$chew = \frac{1}{3} \times 0.48 = 0.16$$

$$cat = \frac{1}{3} \times 0.18 = 0.06$$

cat , cats , fish, dog, chew

DOC1 → [ 0.06, 0.06, 0.06, 0, 0 ] ←

DOC 2 → [ 0 , 0.06, 0.06, 0.06, 0 ]

DOC 3 → [ 0.06, 0, 0, 0.06, 0.16 ]

Advantges

① Intuitive

② fixed size

③ Word Importance is getting Capture

Disadvantages

① sparsity

② fixed vocabulary