# TF-idf

Friday, September 7, 2018        4:01 PM

ypically, the tf-idf weight is composed by two terms: the first computes the normalized Term Frequency (TF), aka. the number of times a word appears in a document, divided by the total number of words in that document; the second term is the Inverse Document Frequency (IDF), computed as the logarithm of the number of the documents in the corpus divided by the number of documents where the specific term appears.

- **TF: Term Frequency**, which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length (aka. the total number of terms in the document) as a way of normalization:

    TF(t) = (Number of times term t appears in a document) / (Total number of terms in the document).

- **IDF: Inverse Document Frequency**, which measures how important a term is. While computing TF, all terms are considered equally important. However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:

    IDF(t) = log_e(Total number of documents / Number of documents with term t in it).

See below for a simple example.

## Example:

Consider a document containing 100 words wherein the word *cat* appears 3 times. The term frequency (i.e., tf) for *cat* is then (3 / 100) = 0.03. Now, assume we have 10 million documents and the word *cat* appears in one thousand of these. Then, the inverse document frequency (i.e., idf) is calculated as log(10,000,000 / 1,000) = 4. Thus, the Tf-idf weight is the product of these quantities: 0.03 * 4 = 0.12.

From <http://www.tfidf.com/>

### Word count

D1: He is a lazy boy. She is also lazy.

D2: Neeraj is a lazy person.

The dictionary created may be a list of unique tokens(words) in the corpus
=['He','She','lazy','boy','Neeraj','person']

Here, D=2, N=6

The count matrix M of size 2 X 6 will be represented as —

|    | He | She | lazy | boy | Neeraj | person |
|----|----|-----|------|-----|--------|--------|
| D1 | 1  | 1   | 2    | 1   | 0      | 0      |
| D2 | 0  | 0   | 1    | 0   | 1      | 1      |

## 2.1.2 TF-IDF vectorization

This is another method which is based on the frequency method **but it is different to the count vectorization in the sense that it takes into account not just the occurrence of a word in a single document but in the entire corpus**. So, what is the rationale behind this? Let us try to understan