

Open in app ↗

Sign up

Sign in



Search

Write



# Text Vectorization: Term Frequency — Inverse Document Frequency (TFIDF)

A technique for converting text into finite length vectors



Vaibhav Jayaswal · [Follow](#)

Published in Towards Data Science · 4 min read · Oct 4, 2020



70



1

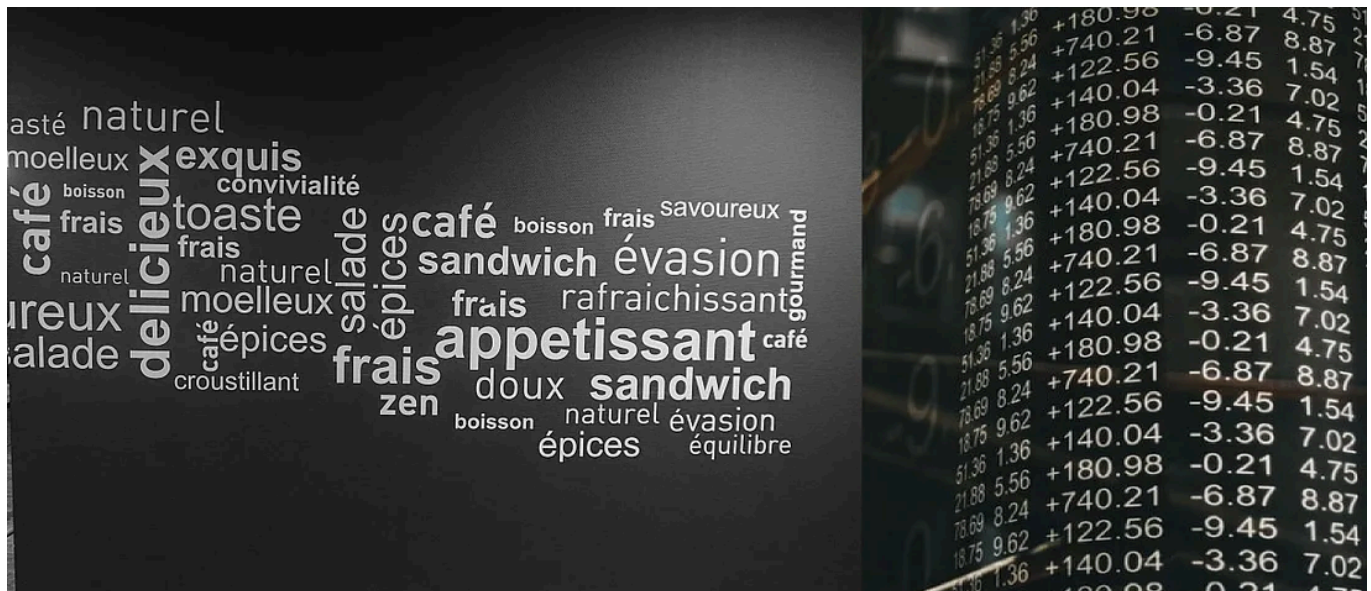


Image by [Hermaion](#) (left) and [Tyler Easton](#) (right)

Bag of words (BoW) converts the text into a feature vector by counting the occurrence of words in a document. It is not considering the importance of words. **Term frequency — Inverse document frequency (TFIDF)** is based on the Bag of Words (BoW) model, which contains insights about the less relevant and more relevant words in a document. The importance of a word in the text is of great significance in information retrieval.

**Example** — If you search something on the search engine, with the help of TFIDF values, search engines can give us the most relevant documents related to our search.

We will be discussing in detail how TFIDF can tell us which word is more important? We will first look into term frequency (TF) and inverse document frequency (IDF) separately and then combine it at the end.

### **Term Frequency (TF)**

It is a measure of the frequency of a word ( $w$ ) in a document ( $d$ ). TF is defined as the ratio of a word's occurrence in a document to the total number of words in a document. The denominator term in the formula is to normalize since all the corpus documents are of different lengths.

$$TF(w, d) = \frac{\text{occurences of } w \text{ in document } d}{\text{total number of words in document } d}$$

Image by Author

### **EXAMPLE**

Documents	Text	Total number of words in a document
A	Jupiter is the largest planet	5
B	Mars is the fourth planet from the sun	8

Image by Author

The initial step is to make a vocabulary of unique words and calculate TF for each document. TF will be more for words that frequently appear in a document and less for rare words in a document.

Words	TF (for A)	TF (for B)
Jupiter	1/5	0
Is	1/5	1/8
The	1/5	2/8
largest	1/5	0
Planet	1/5	1/8
Mars	0	1/8
Fourth	0	1/8
From	0	1/8
Sun	0	1/8

Image by Author

## Inverse Document Frequency (IDF)

It is the measure of the importance of a word. Term frequency (TF) does not consider the importance of words. Some words such as 'of', 'and', etc. can be most frequently present but are of little significance. IDF provides weightage to each word based on its frequency in the corpus D.

IDF of a word (w) is defined as

$$IDF(w, D) = \ln\left(\frac{\text{Total number of documents (N) in corpus D}}{\text{number of documents containing w}}\right)$$

Image by Author

In our example, since we have two documents in the corpus, N=2.

Words	TF (for A)	TF (for B)	IDF
Jupiter	1/5	0	$\ln(2/1) = 0.69$
Is	1/5	1/8	$\ln(2/2) = 0$
The	1/5	2/8	$\ln(2/2) = 0$
largest	1/5	0	$\ln(2/1) = 0.69$
Planet	1/5	1/8	$\ln(2/2) = 0$
Mars	0	1/8	$\ln(2/1) = 0.69$
Fourth	0	1/8	$\ln(2/1) = 0.69$
From	0	1/8	$\ln(2/1) = 0.69$
Sun	0	1/8	$\ln(2/1) = 0.69$

Image by Author

## Term Frequency — Inverse Document Frequency (TFIDF)

It is the product of TF and IDF.

- TFIDF gives more weightage to the word that is rare in the corpus (all the documents).
- TFIDF provides more importance to the word that is more frequent in the document.

$$TFIDF(w, d, D) = TF(w, d) * IDF(w, D)$$

Image by Author

Words	TF (for A)	TF (for B)	IDF	TFIDF (A)	TFIDF (B)
Jupiter	1/5	0	$\ln(2/1) = 0.69$	0.138	0
Is	1/5	1/8	$\ln(2/2) = 0$	0	0
The	1/5	2/8	$\ln(2/2) = 0$	0	0
largest	1/5	0	$\ln(2/1) = 0.69$	0.138	0
Planet	1/5	1/8	$\ln(2/2) = 0$	0.138	0
Mars	0	1/8	$\ln(2/1) = 0.69$	0	0.086
Fourth	0	1/8	$\ln(2/1) = 0.69$	0	0.086
From	0	1/8	$\ln(2/1) = 0.69$	0	0.086
Sun	0	1/8	$\ln(2/1) = 0.69$	0	0.086

Image by Author

After applying TFIDF, text in A and B documents can be represented as a TFIDF vector of dimension equal to the vocabulary words. The value corresponding to each word represents the importance of that word in a particular document.

## Why are we using Ln in the IDF formula?

TFIDF is the product of TF with IDF. Since TF values lie between 0 and 1, not using  $\ln$  can result in high IDF for some words, thereby dominating the TFIDF. We don't want that, and therefore, we use  $\ln$  so that IDF should not completely dominate the TFIDF.

## Disadvantage of TFIDF

It is unable to capture the **semantics**. For example, **funny** and **humorous** are synonyms, but TFIDF does not capture that. Moreover, TFIDF can be computationally expensive if the vocabulary is vast.

## Conclusion

Term Frequency — Inverse Document Frequency (TFIDF) is a technique for text vectorization based on the Bag of words (BoW) model. It performs better than the BoW model as it considers the importance of the word in a document into consideration. The main limitation is that it does not capture the semantic meaning of the words. This limitation of TFIDF can be overcome by more advanced techniques such as word2Vec.

Thanks for reading!

[Tf Idf](#)[Text Vectorization](#)[Term Frequency](#)[NLP](#)