

Bag-of-Words (BoW) and TF-IDF for Creating Features from Text

The Challenge of Making Machines Understand Text

- Machines simply cannot process text data in raw form. They need us to break down the text into a numerical format easily readable by the machine.
- Both **BoW** and **TF-IDF** are techniques that help us convert text sentences into numeric vectors

Here's a sample of reviews about a particular horror movie:

- Review 1: This movie is very scary and long
- Review 2: This movie is not scary and is slow
- Review 3: This movie is spooky and good

Creating Vectors from Text

Word Embedding is one such technique where we can represent the text using vectors. The more popular forms of word embeddings are:

- BoW, which stands for Bag of Words
- TF-IDF, which stands for Term Frequency-Inverse Document Frequency

Bag of Words (BoW) Model

- We will first build a vocabulary from all the unique words in the above three reviews. The vocabulary consists of these 11 words: 'This', 'movie', 'is', 'very', 'scary', 'and', 'long', 'not', 'slow', 'spooky', 'good'.

	1 This	2 movie	3 is	4 very	5 scary	6 and	7 long	8 not	9 slow	10 spooky	11 good	Length of the review(in words)
Review 1	1	1	1	1	1	1	1	0	0	0	0	7
Review 2	1	1	2	0	0	1	1	0	1	0	0	8
Review 3	1	1	1	0	0	0	1	0	0	1	1	6

- Review 1: This movie is very scary and long
- Review 2: This movie is not scary and is slow
- Review 3: This movie is spooky and good

Vector of Review 1: [1 1 1 1 1 1 1 0 0 0 0]

Vector of Review 2: [1 1 2 0 0 1 1 0 1 0 0]

Vector of Review 3: [1 1 1 0 0 0 1 0 0 1 1]

Drawbacks of using a Bag-of-Words (BoW) Model

- 1.If the new sentences contain new words, then our vocabulary size would increase and thereby, the length of the vectors would increase too.
- 2.Additionally, the vectors would also contain many 0s, thereby resulting in a sparse matrix (which is what we would like to avoid)
- 3.We are retaining no information on the grammar of the sentences nor on the ordering of the words in the text.

Term Frequency-Inverse Document Frequency (TF-IDF)

Term frequency–inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.

Term Frequency (TF): It is a measure of how frequently a term, t , appears in a document, d

$$tf_{t,d} = \frac{n_{t,d}}{\text{Number of terms in the document}}$$

n is the number of times the term “ t ” appears in the document “ d ”. Thus, each document and term would have its own TF value.

Review 2: This movie is not scary and is slow

- Vocabulary: ‘This’, ‘movie’, ‘is’, ‘very’, ‘scary’, ‘and’, ‘long’, ‘not’, ‘slow’, ‘spooky’, ‘good’
- Number of words in Review 2 = 8
- TF for the word ‘this’ = (number of times ‘this’ appears in review 2)/(number of terms in review 2) = 1/8

- Review 1: This movie is very scary and long
- Review 2: This movie is not scary and is slow
- Review 3: This movie is spooky and good

Term	Review 1	Review 2	Review 3	TF (Review 1)	TF (Review 2)	TF (Review 3)
This	1	1	1	1/7	1/8	1/6
movie	1	1	1	1/7	1/8	1/6
is	1	2	1	1/7	1/4	1/6
very	1	0	0	1/7	0	0
scary	1	1	0	1/7	1/8	0
and	1	1	1	1/7	1/8	1/6
long	1	0	0	1/7	0	0
not	0	1	0	0	1/8	0
slow	0	1	0	0	1/8	0
spooky	0	0	1	0	0	1/6
good	0	0	1	0	0	1/6

Inverse Document Frequency (IDF)

- IDF is a measure of how important a term is. We need the IDF value because computing just the TF alone is not sufficient to understand the importance of words

$$idf_t = \log \frac{\text{number of documents}}{\text{number of documents with term 't'}}$$

IDF values for all the words in Review 2:

$$IDF('this') = \log(\text{number of documents} / \text{number of documents containing the word 'this'}) = \log(3/3) = \log(1) = 0$$

Similarly,

- $IDF('movie',) = \log(3/3) = 0$
- $IDF('is') = \log(3/3) = 0$
- $IDF('not') = \log(3/1) = \log(3) = 0.48$
- $IDF('scary') = \log(3/2) = 0.18$
- $IDF('and') = \log(3/3) = 0$
- $IDF('slow') = \log(3/1) = 0.48$

- Review 1: This movie is very scary and long
- Review 2: This movie is not scary and is slow
- Review 3: This movie is spooky and good

Term	Review 1	Review 2	Review 3	IDF
This	1	1	1	0.00
movie	1	1	1	0.00
is	1	2	1	0.00
very	1	0	0	0.48
scary	1	1	0	0.18
and	1	1	1	0.00
long	1	0	0	0.48
not	0	1	0	0.48
slow	0	1	0	0.48
spooky	0	0	1	0.48
good	0	0	1	0.48

Hence, we see that words like “is”, “this”, “and”, etc., are reduced to 0 and have little importance; while words like “scary”, “long”, “good”, etc. are words with more importance and thus have a higher value.

We can now compute the TF-IDF score for each word in the corpus. Words with a higher score are more important, and those with a lower score are less important:

We can now compute the TF-IDF score for each word in the corpus. Words with a higher score are more important, and those with a lower score are less important:

$$(tf_idf)_{t,d} = tf_{t,d} * idf_t$$

- Review 1: This movie is very scary and long
- Review 2: This movie is not scary and is slow
- Review 3: This movie is spooky and good

We can now calculate the TF-IDF score for every word in Review 2:

$$TF\text{-}IDF(\text{'this'}, \text{Review 2}) = TF(\text{'this'}, \text{Review 2}) * IDF(\text{'this'}) = 1/8 * 0 = 0$$

Similarly,

- $TF\text{-}IDF(\text{'movie'}, \text{Review 2}) = 1/8 * 0 = 0$
- $TF\text{-}IDF(\text{'is'}, \text{Review 2}) = 1/4 * 0 = 0$
- $TF\text{-}IDF(\text{'not'}, \text{Review 2}) = 1/8 * 0.48 = 0.06$
- $TF\text{-}IDF(\text{'scary'}, \text{Review 2}) = 1/8 * 0.18 = 0.023$
- $TF\text{-}IDF(\text{'and'}, \text{Review 2}) = 1/8 * 0 = 0$
- $TF\text{-}IDF(\text{'slow'}, \text{Review 2}) = 1/8 * 0.48 = 0.06$

Term	Review 1	Review 2	Review 3	IDF	TF-IDF (Review 1)	TF-IDF (Review 2)	TF-IDF (Review 3)
This	1	1	1	0.00	0.000	0.000	0.000
movie	1	1	1	0.00	0.000	0.000	0.000
is	1	2	1	0.00	0.000	0.000	0.000
very	1	0	0	0.48	0.068	0.000	0.000
scary	1	1	0	0.18	0.025	0.022	0.000
and	1	1	1	0.00	0.000	0.000	0.000
long	1	0	0	0.48	0.068	0.000	0.000
not	0	1	0	0.48	0.000	0.060	0.000
slow	0	1	0	0.48	0.000	0.060	0.000
spooky	0	0	1	0.48	0.000	0.000	0.080
good	0	0	1	0.48	0.000	0.000	0.080

Conclusion

Bag of Words just creates a set of vectors containing the count of word occurrences in the document (reviews), while the TF-IDF model contains information on the more important words and the less important ones as well.

Bag of Words vectors are easy to interpret. However, TF-IDF usually performs better in machine-learning models.

- While both Bag-of-Words and TF-IDF have been popular in their own regard, there still remained a void where understanding the context of words was concerned. Detecting the similarity between the words 'spooky' and 'scary', or translating our given documents into another language, requires a lot more information on the documents.
- This is where Word Embedding techniques such as Word2Vec, Continuous Bag of Words (CBOW), Skipgram, etc. come in. You can find a detailed guide to such techniques [here](#):