



Towards Data
Science

Sharing concepts,
ideas, and codes.

Follow



110

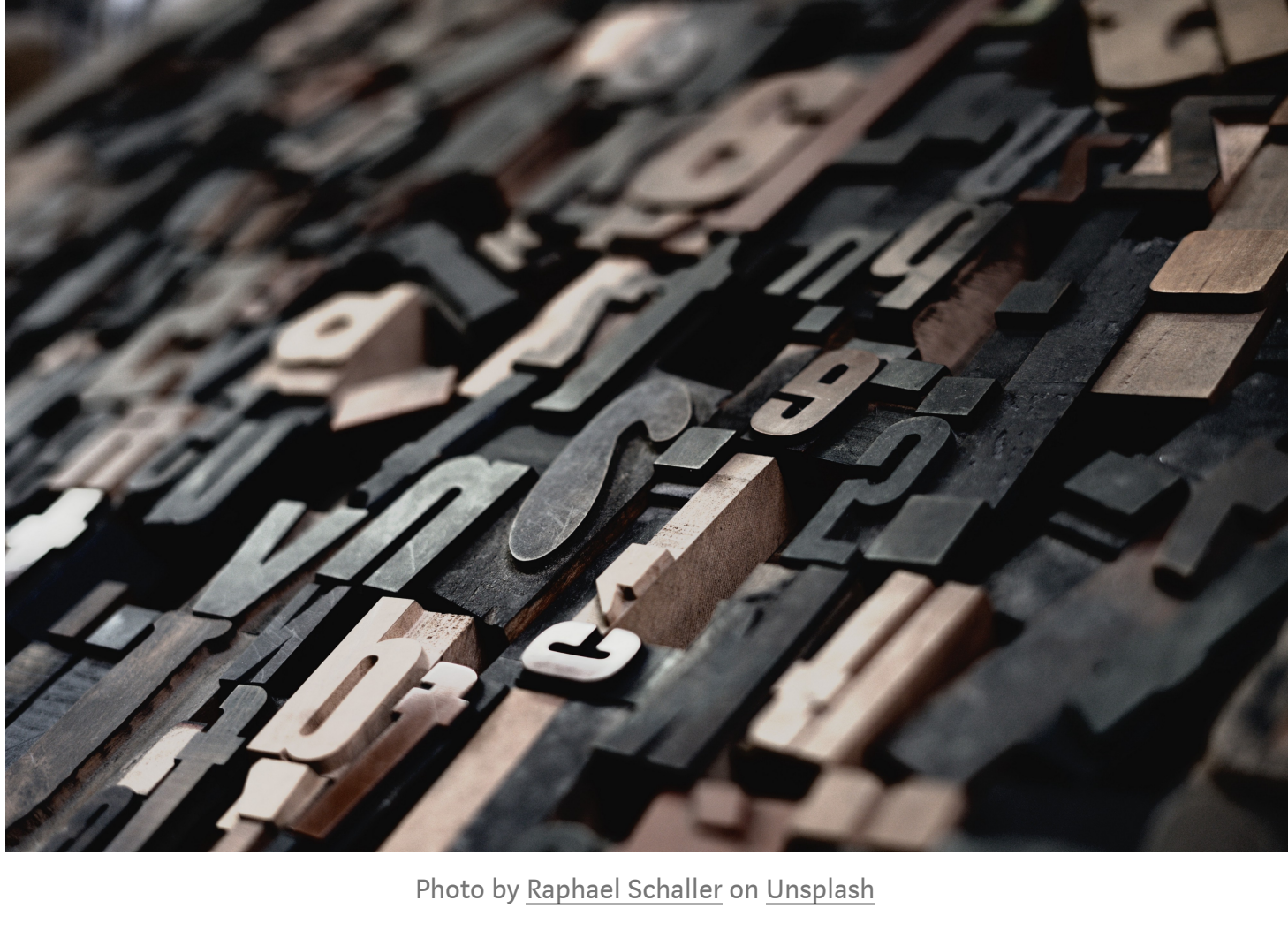


Photo by Raphael Schaller on Unsplash

Get one more story in your member preview when you sign up. It's free.

Sign up with Google

Sign up with Facebook

Already have an account? [Sign in](#)

Natural Language Processing (NLP) is a sub-field of artificial intelligence that deals understanding and processing human language. In light of new advancements in machine learning, many organizations have begun applying natural language processing for translation, chatbots and candidate filtering.

Without further delay let's dive into some code. To start, we'll import the necessary libraries.

```
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
```

In this article, we'll be working with two simple documents containing one sentence each.

```
documentA = 'the man went out for a walk'
documentB = 'the children sat around the fire'
```

Machine learning algorithms cannot work with raw text directly. Rather, the text must be converted into vectors of numbers. In natural language processing, a common technique for extracting features from text is to place all of the words that occur in the text in a bucket. This approach is called a **bag of words** model or **BoW** for short. It's referred to as a *“bag”* of words because any information about the structure of the sentence is lost.

```
bagOfWordsA = documentA.split(' ')
bagOfWordsB = documentB.split(' ')
```

```
{'the', 'man', 'went', 'out', 'for', 'a', 'walk'}
```

By casting the bag of words to a set, we can automatically remove any duplicate words.

```
uniqueWords = set(bagOfWordsA).union(set(bagOfWordsB))
```

Next, we'll create a dictionary of words and their occurrence for each document in the corpus (collection of documents).

```
numOfWordsA = dict.fromkeys(uniqueWords, 0)

for word in bagOfWordsA:
    numOfWordsA[word] += 1

numOfWordsB = dict.fromkeys(uniqueWords, 0)

for word in bagOfWordsB:
    numOfWordsB[word] += 1
```

	a	around	children	fire	for	man	out	sat	the	walk	went
0	1	0	0	0	1	1	1	0	1	1	1
1	0	1	1	1	0	0	0	1	2	0	0

Another problem with the bag of words approach is that it doesn't account for noise. In other words, certain words are used to formulate sentences but do not add any semantic meaning to the text. For example, the most commonly used word in the english language is *the* which represents 7% of all words written or spoken. You couldn't make deduce anything about a text given the fact that it contains the word **the**. On the other hand, words like **good** and **awesome** could be used to determine whether a rating was positive or not.

In natural language processing, useless words are referred to as stop words. The python **natural language toolkit** library provides a list of english stop words.

```
from nltk.corpus import stopwords

stopwords.words('english')
```

{'ourselves', 'hers', 'between', 'yourself', 'but', 'again', 'there', 'about', 'once', 'during', 'out', 'very', 'having', 'with', 'they', 'own', 'an', 'be', 'some', 'for', 'do', 'is', 'yours', 'such', 'into', 'of', 'most', 'itself', 'other', 'off', 'is', 's', 'am', 'or', 'who', 'as', 'toward', 'him', 'each', 'the', 'them', 'selves', 'until', 'below', 'are', 'we', 'these', 'your', 'his', 'through', 'don', 'nor', 'me', 'were', 'her', 'more', 'himself', 'this', 'down', 'should', 'our', 'their', 'while', 'above', 'both', 'up', 'to', 'ours', 'had', 'she', 'all', 'no', 'when', 'at', 'any', 'before', 'them', 'same', 'and', 'been', 'have', 'it', 'will', 'or', 'does', 'yourself', 'their', 'that', 'because', 'what', 'over', 'why', 'so', 'can', 'did', 'not', 'now', 'under', 'be', 'you', 'theseff', 'has', 'just', 'where', 'too', 'only', 'myself', 'which', 'those', 'i', 'after', 'few', 'whom', 'i', 'being', 'if', 'theirs', 'my', 'against', 'a', 'by', 'doing', 'it', 'how', 'further', 'was', 'here', 'than'}

Often times, when building a model with the goal of understanding text, you'll see all of stop words being removed. Another strategy is to score the relative importance of words using TF-IDF.

Term Frequency (TF)

The number of times a word appears in a document divided by the total number of words in the document. Every document has its own term frequency.

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}}$$

The following code implements term frequency in python.

```
def computeTF(wordDict, bagOfWords):
    tfDict = {}
    bagOfWordsCount = len(bagOfWords)
    for word, count in wordDict.items():
        tfDict[word] = count / float(bagOfWordsCount)
    return tfDict
```

The following lines compute the term frequency for each of our documents.

```
tFA = computeTF(numOfWordsA, bagOfWordsA)
tFB = computeTF(numOfWordsB, bagOfWordsB)
```

Inverse Data Frequency (IDF)

The log of the number of documents divided by the number of documents that contain the word **w**. Inverse data frequency determines the weight of rare words across all documents in the corpus.

$$idf(w) = \log\left(\frac{N}{df_t}\right)$$

The following code implements inverse data frequency in python.

```
def computeIDF(documents):
    import math
    N = len(documents)

    idfDict = dict.fromkeys(documents[0].keys(), 0)
    for document in documents:
        for word, val in document.items():
            if val > 0:
                idfDict[word] += 1

    for word, val in idfDict.items():
        idfDict[word] = math.log(N / float(val))
    return idfDict
```

The IDF is computed once for all documents.

```
idfs = computeIDF([numOfWordsA, numOfWordsB])
```

Lastly, the TF-IDF is simply the TF multiplied by IDF.

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

```
def computeTFIDF(tfBagOfWords, idfs):
    tfidf = {}
    for word, val in tfBagOfWords.items():
        tfidf[word] = val * idfs[word]
    return tfidf
```

Finally, we can compute the TF-IDF scores for all the words in the corpus.

```
tfidfA = computeTFIDF(tFA, idfs)
tfidfB = computeTFIDF(tFB, idfs)

df = pd.DataFrame([tfidfA, tfidfB])
```



Rather than manually implementing TF-IDF ourselves, we could use the class provided by sklearn.

```
vectorizer = TfidfVectorizer()

vectors = vectorizer.fit_transform([documentA, documentB])

feature_names = vectorizer.get_feature_names()

dense = vectors.todense()

denselist = dense.tolist()

df = pd.DataFrame(denselist, columns=feature_names)
```



The values differ slightly because sklearn uses a smoothed version idf and various other little optimizations. In an example with more text, the score for the word *the* would be greatly reduced.

Machine Learning | Natural Language Process | TF IDf Python | TF IDf Explained | Tfidf Vectorizer



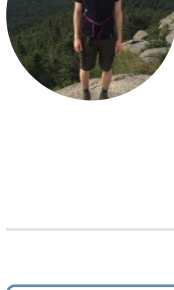
110 claps



See more stories from Towards Data Science.

Create a free Medium account to follow Towards Data Science. You'll see more of their stories on Medium and in your inbox.

Follow



WRITTEN BY

Cory Maklin

Data Science | Data Engineer @ Interset | LinkedIn:
<https://www.linkedin.com/in/cory-maklin>

Follow

See responses (2)

More From Medium

More from Towards Data Science



How To Fake Being a Good Programmer



Sten Sootla in Towar...
Oct 30 · 5 min read ★



3.2K



More from Towards Data Science



The Most Undervalued Standard Python Library



Tyler Folkman in Towar...
Oct 27 · 3 min read ★



3.5K



More from Towards Data Science



Want a data science job? Use the weekend project principle to get it



Daniel Bourke in Towar...
Nov 3 · 4 min read ★



670



Discover Medium

Welcome to a place where words matter. On Medium, smart voices and original ideas take center stage - with no ads in sight. [Watch](#)

Make Medium yours

Follow all the topics you care about, and we'll deliver the best stories for you to your homepage and inbox. [Explore](#)

Become a member

Get unlimited access to the best stories on Medium — and support writers while you're at it. Just \$5/month. [Upgrade](#)