



## Text clustering with K-means and tf-idf



Mihail Salnikov [Follow](#)

Aug 5, 2018 · 4 min read

*First of all, I'm not a native English speaker, then I will probably make a lot of mistakes, sorry about that.*

In this post, I'll try to describe how to clustering text with knowledge, how important word is to a string. Same words in different strings can be badly affected to clustering this kind of data isn't important for deciding. The first part of this publication is the general information about TF-IDF with examples on Python. In the second part, I'll provide you the example showed how this approach can be applied to real tasks.

TF-IDF is useful for clustering tasks, like a document clustering or in other words, tf-idf can help you understand what kind of document you got now.

### TF-IDF

**Term Frequency-Inverse Document Frequency** is a numerical statistic that demonstrates how important a word is to a corpus.

Term Frequency is just ratio number of current word to the number of all words in document/string/etc.

$$tf(t, d) = \frac{n_t}{\sum_k n_k}$$

Frequency of term  $t_i$ , where  $n_t$ —the number of  $t_i$  in current document/string, the sum of  $n_k$  is the number of all terms in current

document/string.

Inverse Document Frequency is a log of the ratio of the number of all documents/string in the corpus to the number of documents with term  $t_i$ .

$$idf(t, D) = \log \frac{|D|}{|\{d_i \in D \mid t \in d_i\}|}$$

tf-idf( $t, d, D$ ) is the product  $tf(t, d)$  to  $idf(t, D)$ .

$$tf-idf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

If you want more theoretic information about TF-IDF I want advice you read publication on [Wikipedia](#) about it or read [NLP Stanford post](#).

Well, now time for a real example on Python.

## TF-IDF example on Python

For all code below you need [python 3.5 or newer](#) and [scikit-learn](#) and [pandas](#) packages.

Firstly, let's talk about a data set. For this really simple example, I just set a simple corpus with 3 strings. In this example, strings play a role documents.

```
1 corpus = """
2 Simple example with Cats and Mouse
3 Another simple example with dogs and cats
4 Another simple example with mouse and cheese
... """
```

After that lets make [bags of words](#) for our corpus and for every string too. But before we have to clear the data.

```

1  # clearing and tokenizing
2  l_A = corpus[0].lower().split()
3  l_B = corpus[1].lower().split()
4  l_C = corpus[2].lower().split()
5
6  # Calculating bag of words
7  word_set = set(l_A).union(set(l_B)).union(set(l_C))
8
9  word_dict_A = dict.fromkeys(word_set, 0)
10 word_dict_B = dict.fromkeys(word_set, 0)
11 word_dict_C = dict.fromkeys(word_set, 0)
12
13 for word in l_A:
14     word_dict_A[word] += 1

```

	and	another	cats	cheese	dogs	example	mouse	simple	with
0	1	0	1	0	0	1	1	1	1
1	1	1	1	0	1	1	0	1	1
2	1	1	0	1	0	1	1	1	1

that's what we get

In the case of the term frequency, the simplest choice is to use the raw count of a term in a string. For calculating tf for all terms, we must fill a dictionary as follows.

```

1  def compute_tf(word_dict, l):
2      tf = {}
3      sum_nk = len(l)
4      for word, count in word_dict.items():
5          tf[word] = count/sum_nk
6      return tf
7
8

```

idf is a measure of how much information the token or word in our case, provides. For calculating idf we need fill dict too.

```

1  def compute_idf(strings_list):
2      n = len(strings_list)
3      idf = dict.fromkeys(strings_list[0].keys(), 0)
4      for l in strings_list:
5          for word, count in l.items():
6              if count > 0:
7                  idf[word] += 1
8
9      for word, v in idf.items():

```

Now, I remain you that tf-idf is the product of tf to idf. For our python example, tf-idf it dict with the corresponding products.

```

1  def compute_tf_idf(tf, idf):
2      tf_idf = dict.fromkeys(tf.keys(), 0)
3      for word, v in tf.items():
4          tf_idf[word] = v * idf[word]
5      return tf_idf
6
7  tf_idf_A = compute_tf_idf(tf_A, idf)

```

OK, now we have tf-idf weights for each word in our corpus. Below you can clearly see the difference between the original bag of words and the new bag of words with tf-idf weights. For example ‘dogs’, ‘cats’ and ‘mouse’ is important words, but word ‘and’ is not important, because this word is in all the strings and we can’t understand what is a string by the word ‘and’.

	and	another	cats	cheese	dogs	example	mouse	simple	with
0	0.0	0.000000	0.067578	0.000000	0.000000	0.0	0.067578	0.0	0.0
1	0.0	0.057924	0.057924	0.000000	0.156945	0.0	0.000000	0.0	0.0
2	0.0	0.057924	0.000000	0.156945	0.000000	0.0	0.057924	0.0	0.0

TF-IDF bag of words

	and	another	cats	cheese	dogs	example	mouse	simple	with
0	1	0	1	0	0	1	1	1	1
1	1	1	1	0	1	1	0	1	1
2	1	1	0	1	0	1	1	1	1

original bag of words

. . .

## KMeans clustering with TF-IDF weights

Now, when we understand how TF-IDF work the time has come for almost real example of clustering with TF-IDF weights. For real life we can use scikit-learn implementation of TF-IDF and KMeans and I suggest you use implementations from scikit-learn or from another popular libraries or frameworks because it's reducing a number of potential errors in your code.

For this example, we must import TF-IDF and KMeans, added corpus of text for clustering and process its corpus.

```

1  from sklearn.feature_extraction.text import TfidfVectorizer
2  from sklearn.cluster import KMeans
3
4
5  all_textall_text == """
6  Google and Facebook are strangling the free press to death
7  Your 60-second guide to security stuff Google touted today
8  A Guide to Using Android Without Selling Your Soul to Google
9  Review: Lenovo's Google Smart Display is pretty and intelligent
10 Google Maps user spots mysterious object submerged off the
11 Android is better than iOS
12 In information retrieval, tf-idf or TFIDF, short for term frequency
13 is a numerical statistic that is intended to reflect
14 how important a word is to a document in a collection or corpus.
15 It is often used as a weighting factor in searches of information
16 text mining, and user modeling. The tf-idf value increases
17 to the number of times a word appears in the document

```

After that let's fit Tfidf and let's fit KMeans, with scikit-learn it's really easy.

```

1  tfidf_vectorizer = TfidfVectorizer(preprocessor=preprocessin
2  tfidf = tfidf_vectorizer.fit_transform(all_text)
3
4  kmeans = KMeans(n_clusters=2).fit(tfidf)

```

Now we have learned KMeans model with  $k = 2$  for clustering strings, it's easy, right?

For predicting, just use predict method as follows.

```
1 lines_for_predicting = ["tf and idf is awesome!", "some andr  
2 kmeans.predict(tfidf_vectorizer.transform(lines_for_predicti  
3  
4 # array([0. 11. dtype=int32)
```

There we can see, that string *'tf and idf is awesome!'* and *'some androids is there'* from different clusters and it's right.

. . .

In addition, you can read [Jupyter notebook](#) with this examples.

Thanks for the reading, please leave a feedback. This can help me improve the quality of my future posts.

And don't forget to follow me on [twitter](#).

