

KDnuggets

[Subscribe to KDnuggets News](#)



- [Blog/News](#)
- [Opinions](#)
- [Tutorials](#)
- [Top stories](#)
- [Companies](#)
- [Courses](#)
- [Datasets](#)
- [Education](#)
- [Events \(online\)](#)
- [Jobs](#)
- [Software](#)
- [Webinars](#)



Free Online Statistics Course
Build practical skills in using data to solve problems.

[Enroll Now >](#)

[Free Online Statistics Course. Build practical skills in using data to solve problems](#)

NEW Topics: [Coronavirus](#) | [AI](#) | [Data Science](#) | [Deep Learning](#) | [Machine Learning](#) | [Python](#) | [R](#) | [Statistics](#)

[KDnuggets Home](#) » [News](#) » [2018](#) » [Aug](#) » [Tutorials, Overviews](#) » WTF is TF-IDF? ([18:n30](#))

WTF is TF-IDF?

[<= Previous post](#)

[Next post =>](#)

Like 18

Share 18

Tweet

Share

Share

31

Tags: [Information Retrieval](#), [Python](#), [Text Analytics](#), [Text Mining](#), [TF-IDF](#)

Relevant words are not necessarily the most frequent words since stopwords like “the”, “of” or “a” tend to occur very often in many documents.



John Snow LABS
Spark NLP
for Healthcare Data Scientists
HALF-DAY CERTIFICATION COURSES
Save 20% with code 'KDnuggets20'

[SPARK NLP for Healthcare Data Scientists. Half-day certification courses](#)
[Save 20% w. code KDnuggets20](#)

[comments](#)

By [Enrique Fueyo](#), CTO & Co-founder @ Lang.ai

<https://www.kdnuggets.com/2018/08/wtf-tf-idf.html>



Frame from "The Incredibles" (2004) movie

TF-IDF, which stands for **term frequency—inverse document frequency**, is a scoring measure widely used in information retrieval (IR) or summarization. TF-IDF is intended to reflect **how relevant a term is in a given document**.

The intuition behind it is that if a word occurs *multiple times in a document*, we should boost its relevance as it should be more meaningful than other words that appear fewer times (TF). At the same time, if a word occurs many times in a document but also *along many other documents*, maybe it is because this word is just a frequent word; not because it was relevant or meaningful (IDF).

Defining what a “relevant word” means

We can come up with a more or less subjective definition driven by our intuition: a word’s relevance is proportional to the amount of information that it gives about its context (a sentence, a document or a full dataset). That is, the most relevant words are those that would help us, as humans, to better understand a whole document without reading it all.

As pointed out, **relevant words are not necessarily the most frequent words** since stopwords like “the”, “of” or “a” tend to occur very often in many documents.

There is another caveat: if we want to summarize a document compared to a whole dataset about an specific topic (let’s say, movie reviews), there will be words (other than stopwords, like *character* or *plot*), that could occur many times in the document as well as in many other documents. These words are not useful to summarize a document because they convey little discriminating power; they say very little about what the document contains compared to the other documents.

Let’s go through some examples to better illustrate how TF-IDF works.

Search engine example

Let’s suppose we have a database with thousands of cats descriptions and **a user wants to search for furry cats**, so she/he issues the query “*the furry cat*”. As a search engine, we have to decide which documents should be returned from our database.

If we have documents that match the exact query, there is no doubt but... what if we have to decide between partial matches? For simplicity, let’s say we have to choose between these two descriptions:

1. “the lovely cat”
2. “a furry kitten”

The first description contains **2 out of 3 words from the query** and the second one matches **just 1 out of 3**, then we would pick the first description. How can TF-IDF help us to choose the second description instead of the first one?

The TF is the same for each word, no difference here. However, we could expect that **the terms “cat” and “kitten” would appear in many documents** (large document frequency implies low IDF), while **the term “furry” will appear in fewer documents** (larger IDF). So the TF-IDF for cat & kitten has a low value whereas the TF-IDF is larger for “furry”, i.e. in our database **the word “furry” has more discriminative power than “cat” or “kitten”**.

Conclusion

If we use the TF-IDF to weight the different words that matched the query, “furry” would be more relevant than “cat” and so we could eventually choose “the furry kitten” as the best match.

Summarization example

Now, imagine we would like to **automatically summarize some movies using their wikipedia pages**. We plan to create a tagcloud that will help us to understand what each movie is about, so our task is to decide what words to put in it (and their sizes).

Given that a month ago my family adopted a new dog called *Mawi* (named after Moana's character Maui), we will use [Moana's wikipedia page](#) for this example.



My dogs: Cala & Mawi

The first idea would be to make a list with the most frequent words:

Moana

Most frequent: ['film', 'moana', 'the', 'million', 'disney', 'maui', 'day', 'release', 'te', 'animation', 'weekend', 'heart', 'ocean', 'it', 'story', 'island', 'fiti', 'version', 'in', 'animate']

We can see that some words could be useful to grasp a little bit of the film's plot: moana, maui, ocean, island... **but they still provide very little information.** Moreover, they are mixed with other words that are irrelevant and just add noise (like film, million, day, release or weekend).

If we run the same experiment with other movies, we can see that the results are somehow similar regarding irrelevant and noisy words. Wikipedia pages include information that is shared or have a similar structure among all them: They are *animated films* that made *millions* of dollars since the first *day* of their *release*... but this is something we already knew or, at least, it is not relevant for us in order to understand the plot.

Overcoming the problem

To eliminate what is shared among all movies and extract **what individually identifies each one**, TF-IDF should be a very handy tool. With the most frequent words (TF) we got a first approximation, but IDF should help us to refine the previous list and get better results. Taking into account the document frequency, we can try to penalize generic words (these that appeared in many wikipedia pages) reducing their relevance.

Moana

TF-IDF: ['maui', 'te', 'moana', 'fiti', 'cravalho', 'goddess', 'tui', 'polynesian', 'tala', 'kā', 'māori', 'auli'i', 'clement', 'fishhook', 'tamatoa', 'jemaine', 'tattoo', 'dubbing', 'musker', 'clements']

Thanks to TF-IDF words like *film*, *million* or *release* have disappeared from the top of the list and we got some new more meaningful words like *polynesian* or *tattoo*.

We can also run the same analysis with the other reviews and get the following results where, again, the generic words have disappeared in favor of more movie-specific words:



MOANA

Most relevant	TF - IDF
film	maui
moana	te
the	moana
million	fiti
disney	cravalho
maui	goddess
day	tui
release	polynesian



THE INCREDIBLES

Most relevant	TF - IDF
film	parrs
the	syndrome
incredibles	violet
bird	omnidroid
pixar	parr
release	mirage
bob	nomanisan
jack	helen



MONSTERS INC.

Most relevant	TF - IDF
film	sulley
sulley	waternoose
monsters	boo
the	cda
mike	randall
monster	scarer
pixar	fizt
story	celia

Most relevant VS TF-IDF comparison

For reference, these are the sorted lists with more words:

Moana

Most frequent: ['film', 'moana', 'the', 'million', 'disney', 'maui', 'day', 'release', 'te', 'animation', 'weekend', 'heart', 'ocean', 'it', 'story', 'island', 'fiti', 'version', 'in', 'animate']

TF-IDF: ['maui', 'te', 'moana', 'fiti', 'cravalho', 'goddess', 'tui', 'polynesian', 'tala', 'kā', 'māori', 'auli'i', 'clement', 'fishhook', 'tamatoa', 'jemaine', 'tattoo', 'dubbing', 'musker', 'clements']

The Incredibles

Most frequent: ['film', 'the', 'incredibles', 'bird', 'pixar', 'release', 'bob', 'jack', 'i', 'award', 'good', 'animate', 'feature', 'it', 'syndrome', 'character', 'family', 'work', 'superhero', 'animation']

TF-IDF: ['parrs', 'syndrome', 'violet', 'omnidroid', 'parr', 'mirage', 'nomanisan', 'helen', 'bird', 'edna', 'superhero', 'frozone', 'underminer', 'iron', 'wallin', 'suburban', 'metroville', 'dash', 'incredibles', 'incredible']

Monsters, Inc

Most frequent: ['film', 'sulley', 'monsters', 'the', 'mike', 'monster', 'pixar', 'story', 'child', 'randall', 'boo', 'disney', 'character', 'work', 'good', 'in', 'release', 'he', 'docter', 'fur']

TF-IDF: ['sulley', 'waternoose', 'boo', 'cda', 'randall', 'scarer', 'fizt', 'celia', 'kahrs', 'sullivan', 'fur', 'factory', 'scare', 'monster', 'wazowski', 'mike', 'tentacle', 'madrid', 'pidgeon', 'laughter']

How we did it

Here it is the code we have used to obtain the relevant words lists used in this article:

```

1  # Install dependencies first
2  #     pip install spacy wikipedia
3  #     python -m spacy download en
4  # Run
5  #     python tfidf.py
6
7  from collections import Counter
8  import math
9
10 import wikipedia
11 import spacy
12
13 nlp = spacy.load('en')
14
15 pages = [
16     "The Tigger Movie 2000", "Dinosaur 2000 Movie", "The Emperor's New Groove 2000", "Recess School's Out", "Atlantis: The Lost Empire", "Monste
17     "A Christmas Carol 2009", "The Princess and the Frog 2009", "Toy Story", "Toy Story 2", "Toy Story 3", "Tangled", "Mars Needs Moms", "Cars 2
18     "Shrek", "Shrek 2", "Shrek 3", "Antz", "A bugs life", "Bee movie", "Madagascar 2005 film", "Madagascar 2", "Kung fu panda", "Kung fu panda 2
19 ]
20
21 def valid_token(tk):
22     is_valid = tk.is_alpha
23     return is_valid and not tk.is_stop
24
25 def get_lemma(tk):
26     if tk.pos_ == 'PRON' or tk.lemma_ == '-PRON-':
27         return tk.text.lower()
28     return tk.lemma_.lower()
29
30 def read_wikipedia_page(page_name):
31     page = wikipedia.page(page_name)
32     content = page.content
33     return content
34
35 def tokenize_page(page_name):
36     text = read_wikipedia_page(page_name)
37     return [
38         get_lemma(t)
39         for t in nlp(text)
40         if valid_token(t)
41     ]
42
43 vocabulary = set()
44 idf_counter = Counter()
45
46 for page in pages:
47     print("    Processing page {}".format(page))
48     page_words = set(tokenize_page(page))
49     vocabulary = vocabulary | page_words
50     idf_counter.update(page_words)
51
52 print("All pages processed")
53
54 idf = {
55     word: math.log(len(pages)/df, 2)
56     for word, df in idf_counter.items()
57 }
58
59 print("vocabulary size: {}".format(len(vocabulary)))
60
61
62 def analyze_page(target_page):
63     target_words = tokenize_page(target_page)

```

```

64 tfidf = {
65     word: (1 + math.log(_tf, 2)) * idf[word]
66     for word, _tf in Counter(target_words).items()
67 }
68 num_words = 20
69 most_frequent = [
70     w for (w, _) in Counter(target_words).most_common(num_words)
71 ]
72 sorted_tfidf = [
73     w for (w, _) in sorted(tfidf.items(), key=lambda kv: kv[1], reverse=True)
74 ]
75 print(target_page)
76 print("Most frequent: {}".format(most_frequent))
77 print("Higher TF-IDF: {}".format(sorted_tfidf[:num_words]))
78
79
80 analyze_page("Moana")
81 analyze_page("The Incredibles 2004")
82 analyze_page("Monsters, Inc")

```

tfidf.py hosted with ❤ by GitHub

[view raw](#)

Next steps

For the base set of documents we have used a collection of 65 hand-picked movies. This was only for demo purposes but, to get better results with the TF-IDF, we should use a much larger document base.

Some other experiments we could easily try while using this code as a starting point are:

- Try different weighting scores for either term-frequency or document-frequency. How do they affect the list?
- Use a different document base. What would have happened if we had chosen any movies instead of focusing on a single genre?
- How are the results if we try a different domain like music bands?

Check the other articles in our [Building Lang.ai](#) publication. We write about [Machine Learning](#), [Software Development](#), and our [Company Culture](#).

Bio: [Enrique Fueyo](#) is CTO and co-founder at Lang.ai, working on unsupervised AI for language understanding to build products and services that help both companies and developers who have to deal with unstructured text data.

[Original](#). Reposted with permission.

Related:

- [Understanding What is Behind Sentiment Analysis – Part 1](#)
- [Understanding What is Behind Sentiment Analysis – Part 2](#)
- [A General Approach to Preprocessing Text Data](#)

[<= Previous post](#)

[Next post =>](#)

Top Stories Past 30 Days

Most Popular

- [Five Cool Python Libraries for Data Science](#)
- [The Super Duper NLP Repo: 100 Ready-to-Run Colab Notebooks](#)
- [Natural Language Processing Recipes: Best Practices and Examples](#)
- [24 Best \(and Free\) Books To Understand Machine Learning](#)
- [Free High-Quality Machine Learning & Data Science Books & Courses: Quarantine Edition](#)
- [Mathematics for Machine Learning: The Free eBook](#)

Most Shared

- [Free High-Quality Machine Learning & Data Science Books & Courses: Quarantine Edition](#)
- [Beginners Learning Path for Machine Learning](#)
- [Should Data Scientists Model COVID19 and other Biological Events](#)
- [The Super Duper NLP Repo: 100 Ready-to-Run Colab Notebooks](#)
- [Natural Language Processing Recipes: Best Practices and Examples](#)
- [AI and Machine Learning for Healthcare](#)