

Pyevolve

by Christian S. Perone

[Project](#)

[About](#)



Machine Learning :: Text feature extraction (tf-idf) – Part II

🕒 03/10/2011 📁 Machine Learning, Python 🔖 feature extraction, inverse document frequency, machine learning, Python, scikit.learn, sklearn, term frequency, text mining, tf-idf

Read the first part of this tutorial: [Text feature extraction \(tf-idf\) – Part I](#).

This post is a **continuation** of the first part where we started to learn the theory and practice about text feature extraction and vector space model representation. I really recommend you **to read the first part** of the post series in order to follow this second post.

About Me



Recent Comments

jpf on [Machine Learning :: Text feature extraction \(tf-idf\) – Part I](#)

Christian S. Perone on [Machine Learning :: Cosine Similarity for Vector Space Models \(Part III\)](#)

Isura Nirmal on [Machine Learning :: Cosine Similarity for Vector Space Models \(Part III\)](#)

Since a lot of people liked the first part of this tutorial, this second part is a little longer than the first.

Introduction

In the first post, we learned how to use the **term-frequency** to represent textual information in the vector space. However, the main problem with the term-frequency approach is that it scales up frequent terms and scales down rare terms which are empirically more informative than the high frequency terms. The basic intuition is that a term that occurs frequently in many documents is not a good discriminator, and really makes sense (at least in many experimental tests); the important question here is: why would you, in a classification problem for instance, emphasize a term which is almost present in the entire corpus of your documents ?

The tf-idf weight comes to solve this problem. What tf-idf gives is how important is a word to a document in a collection, and that's why tf-idf incorporates local and global parameters, because it takes in consideration not only the isolated term but also the term within the document collection. What tf-idf then does to solve that problem, is to scale down the frequent terms while scaling up the rare terms; a term that occurs 10 times more than another isn't 10 times more important than it, that's why tf-idf uses the logarithmic scale to do that.

But let's go back to our definition of the $tf(t, d)$ which is actually the term count of the term t in the document d . The use of this simple term frequency could lead us to problems like *keyword spamming*, which is when we have a repeated term in a document with the purpose of improving its ranking on an IR (*Information Retrieval*) system or even create a bias towards long documents, making them look more important than they are just because of the high frequency of the term in the document.

To overcome this problem, the term frequency $tf(t, d)$ of a document on a vector space is

Antonio BIHL Vaz > Curitiba on

Rastreamento em tempo real de avioes em Porto Alegre utilizando Raspberry Pi + Radio UHF (SDR RTL2832U)

Antonio on Rastreamento em tempo real de avioes em Porto Alegre utilizando Raspberry Pi + Radio UHF (SDR RTL2832U)

Recent Posts

Real time Drone object tracking using Python and OpenCV

Arduino and OLED display to monitor Redis for fun and profit

Recebendo dados de baloes meteorologicos da Aeronautica

Simple and effective coin segmentation using Python and OpenCV

Despesas de Custeio e Lei de Benford

Universality, primes and space communication

The beauty of Bitcoin P2P network

Protocolcoin – a pure Python Bitcoin protocol implementation

Mapa de calor dos dados de acidentes de transito do DataPoa

Book Suggestion: Codex Seraphinianus

Machine Learning - Cosine Similarity for

usually also normalized. Let's see how we normalize this vector.

Vector normalization

Suppose we are going to normalize the term-frequency vector \vec{v}_{d4} that we have calculated in the first part of this tutorial. The document $d4$ from the first part of this tutorial had this textual representation:

1. d4: We can see the shining sun, the bright sun.

And the vector space representation using the non-normalized term-frequency of that document was:

$$\vec{v}_{d4} = (0, 2, 1, 0)$$

To normalize the vector, is the same as calculating the **Unit Vector** of the vector, and they are denoted using the “hat” notation: \hat{v} . The definition of the unit vector \hat{v} of a vector \vec{v} is:

$$\hat{v} = \frac{\vec{v}}{\|\vec{v}\|_p}$$

Where the \hat{v} is the unit vector, or the normalized vector, the \vec{v} is the vector going to be normalized and the $\|\vec{v}\|_p$ is the norm (magnitude, length) of the vector \vec{v} in the L^p space (don't worry, I'm going to explain it all).

The unit vector is actually nothing more than a normalized version of the vector, is a vector which the length is 1.

[Machine Learning :: Cosine Similarity for Vector Space Models \(Part III\)](#)

[Rastreamento em tempo real de avioes em Porto Alegre utilizando Raspberry Pi + Radio UHF \(SDR RTL2832U\)](#)

[Raspberry Pi & Arduino: a laser pointer communication and a LDR voltage sigmoid](#)

[Machine Learning :: Cosine Similarity for Vector Space Models \(Part III\)](#)

2 comments

[Machine Learning :: Text feature extraction \(tf-idf\) – Part I](#)

2 comments

[Real time Drone object tracking using Python and OpenCV](#)

0 comments

[Machine Learning :: Text feature extraction \(tf-idf\) – Part II](#)

1 comment

[Simple and effective coin segmentation using Python and OpenCV](#)

1 comment

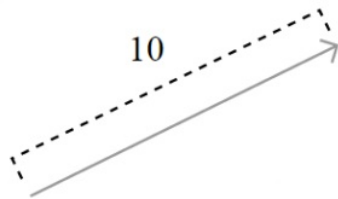
[Raspberry Pi & Arduino: a laser pointer communication and a LDR voltage sigmoid](#)

0 comments

[Genetic Programming meets Python](#)

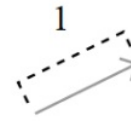
0 comments

This vector \vec{v}
has a length of 10.



\Rightarrow
Process of
Normalization

This vector \vec{v}
points in the same
direction and now
has a length of 1.
(unit vector)



*The normalization process (Source:
<http://processing.org/learning/pvector/>)*

But the important question here is how the length of the vector is calculated and to understand this, you must understand the motivation of the L^p spaces, also called **Lebesgue spaces**.

Lebesgue spaces

Rastreamento em tempo real de avioes em
Porto Alegre utilizando Raspberry Pi +
Radio UHF (SDR RTL2832U)

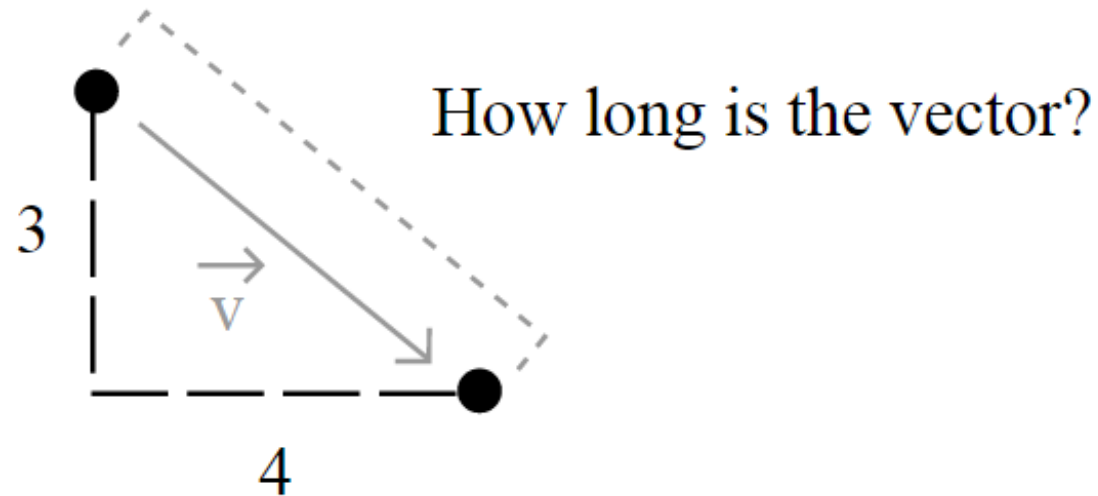
2 comments

Hacking into Python objects internals

1 comment

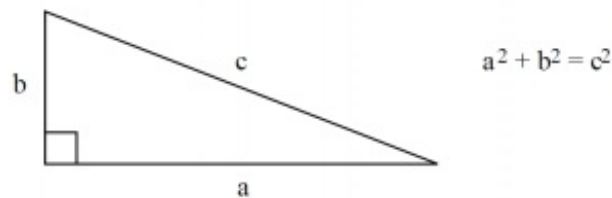
Python: acessing near real-time MODIS
images and fire data from NASA's Aqua and
Terra satellites

0 comments



How long is this vector ? (Source: <http://processing.org/learning/pvector/>)

Usually, the length of a vector $\vec{u} = (u_1, u_2, u_3, \dots, u_n)$ is calculated using the **Euclidean norm** – a norm is a function that assigns a strictly positive length or size to all vectors in a vector space -, which is defined by:



(Source: <http://processing.org/learning/pvector/>)

$$\|\vec{u}\| = \sqrt{u_1^2 + u_2^2 + u_3^2 + \dots + u_n^2}$$

But this isn't the only way to define length, and that's why you see (sometimes) a number p together with the norm notation, like in $\|\vec{u}\|_p$. That's because it could be generalized as:

$$\|\vec{u}\|_p = (|u_1|^p + |u_2|^p + |u_3|^p + \dots + |u_n|^p)^{\frac{1}{p}}$$

and simplified as:

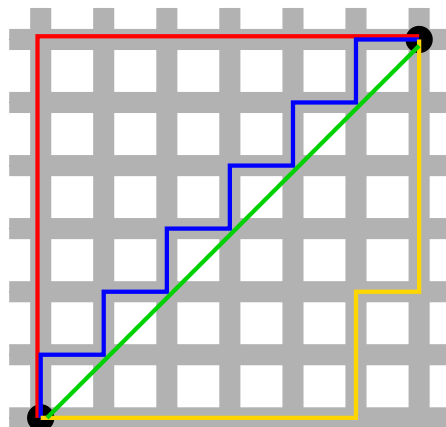
$$\|\vec{u}\|_p = \left(\sum_{i=1}^n |\vec{u}_i|^p \right)^{\frac{1}{p}}$$

So when you read about a **L2-norm**, you're reading about the **Euclidean norm**, a norm with $p = 2$, the most common norm used to measure the length of a vector, typically called "magnitude"; actually, when you have an unqualified length measure (without the p number), you have the **L2-norm** (Euclidean norm).

When you read about a **L1-norm**, you're reading about the norm with $p = 1$, defined as:

$$\|\vec{u}\|_1 = (|u_1| + |u_2| + |u_3| + \dots + |u_n|)$$

Which is nothing more than a simple sum of the components of the vector, also known as **Taxicab distance**, also called Manhattan distance.





Taxicab geometry versus Euclidean distance: In taxicab geometry all three pictured lines have the same length (12) for the same route. In Euclidean geometry, the green line has length $6 \times \sqrt{2} \approx 8.48$, and is the unique shortest path.

Source: [Wikipedia :: Taxicab Geometry](#)

Note that you can also use any norm to normalize the vector, but we're going to use the most common norm, the L2-Norm, which is also the default in the 0.9 release of the [scikits.learn](#). You can also find papers comparing the performance of the two approaches among other methods to normalize the document vector, actually you can use any other method, but you have to be concise, once you've used a norm, you have to use it for the whole process directly involving the norm (*a unit vector that used a L1-norm isn't going to have the length 1 if you're going to take its L2-norm later*).

Back to vector normalization

Now that you know what the vector normalization process is, we can try a concrete example, the process of using the L2-norm (we'll use the right terms now) to normalize our vector $\vec{v}_{d_4} = (0, 2, 1, 0)$ in order to get its unit vector \hat{v}_{d_4} . To do that, we'll simply plug it into the definition of the unit vector to evaluate it:

$$\hat{v} = \frac{\vec{v}}{\|\vec{v}\|_p}$$

$$\hat{v}_{d_4} = \frac{\vec{v}_{d_4}}{\|\vec{v}_{d_4}\|_2}$$

$$\hat{v}_{d_4} = \frac{(0,2,1,0)}{\sqrt{0^2+2^2+1^2+0^2}}$$

$$\hat{v}_{d_4} = \frac{(0,2,1,0)}{\sqrt{5}}$$

$$\hat{v}_{d_4} = (0.0, 0.89442719, 0.4472136, 0.0)$$

And that is it! Our normalized vector \hat{v}_{d_4} has now a L2-norm $\|\hat{v}_{d_4}\|_2 = 1.0$.

Note that here we have normalized our term frequency document vector, but later we're going to do that after the calculation of the tf-idf.

The term frequency - inverse document frequency (tf-idf) weight

Now you have understood how the vector normalization works in theory and practice, let's continue our tutorial. Suppose you have the following documents in your collection (taken from the first part of tutorial):

1. Train Document Set:
- 2.
3. d1: The sky is blue.
4. d2: The sun is bright.
- 5.
6. Test Document Set:
- 7.
8. d3: The sun in the sky is bright.
9. d4: We can see the shining sun, the bright sun.

Your document space can be defined then as $D = \{d_1, d_2, \dots, d_n\}$ where n is the number of documents in your corpus, and in our case as $D_{train} = \{d_1, d_2\}$ and $D_{test} = \{d_3, d_4\}$. The cardinality of our document space is defined by $|D_{train}| = 2$ and $|D_{test}| = 2$, since we have only 2 two documents for training and testing, but they obviously don't need to have the same cardinality.

Let's see now, how idf (inverse document frequency) is then defined:

$$\text{idf}(t) = \log \frac{|D|}{1 + |\{d : t \in d\}|}$$

where $|\{d : t \in d\}|$ is the **number of documents** where the term t appears, when the term-frequency function satisfies $\text{tf}(t, d) \neq 0$, we're only adding 1 into the formula to avoid zero-division.

The formula for the tf-idf is then:

$$\text{tf-idf}(t) = \text{tf}(t, d) \times \text{idf}(t)$$

and this formula has an important consequence: a high weight of the tf-idf calculation is reached when you have a high term frequency (tf) in the given document (*local parameter*) and a low document frequency of the term in the whole collection (*global parameter*).

Now let's calculate the idf for each feature present in the feature matrix with the term frequency we have calculated in the first tutorial:

$$M_{train} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 2 & 1 & 0 \end{bmatrix}$$

Since we have 4 features, we have to calculate $\text{idf}(t_1), \text{idf}(t_2), \text{idf}(t_3), \text{idf}(t_4)$:

$$\text{idf}(t_1) = \log \frac{|D|}{1+|\{d:t_1 \in d\}|} = \log \frac{2}{1} = 0.69314718$$

$$\text{idf}(t_2) = \log \frac{|D|}{1+|\{d:t_2 \in d\}|} = \log \frac{2}{3} = -0.40546511$$

$$\text{idf}(t_3) = \log \frac{|D|}{1+|\{d:t_3 \in d\}|} = \log \frac{2}{3} = -0.40546511$$

$$\text{idf}(t_4) = \log \frac{|D|}{1+|\{d:t_4 \in d\}|} = \log \frac{2}{2} = 0.0$$

These idf weights can be represented by a vector as:

$$\vec{idf}_{train} = (0.69314718, -0.40546511, -0.40546511, 0.0)$$

Now that we have our matrix with the term frequency (M_{train}) and the vector representing the idf for each feature of our matrix (\vec{idf}_{train}), we can calculate our tf-idf weights. What we have to do is a simple multiplication of each column of the matrix M_{train} with the respective \vec{idf}_{train} vector dimension. To do that, we can create a square **diagonal matrix** called M_{idf} with both the vertical and horizontal dimensions equal to the vector \vec{idf}_{train} dimension:

$$M_{idf} = \begin{bmatrix} 0.69314718 & 0 & 0 & 0 \\ 0 & -0.40546511 & 0 & 0 \\ 0 & 0 & -0.40546511 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

and then multiply it to the term frequency matrix, so the final result can be defined then as:

$$M_{tf-idf} = M_{train} \times M_{idf}$$

Please note that the matrix multiplication isn't commutative, the result of $A \times B$ will be different than the result of the $B \times A$, and this is why the M_{idf} is on the right side of the multiplication, to accomplish the desired effect of multiplying each idf value to its corresponding feature:

$$\begin{bmatrix} tf(t_1, d_1) & tf(t_2, d_1) & tf(t_3, d_1) & tf(t_4, d_1) \\ tf(t_1, d_2) & tf(t_2, d_2) & tf(t_3, d_2) & tf(t_4, d_2) \end{bmatrix} \times \begin{bmatrix} idf(t_1) & 0 & 0 & 0 \\ 0 & idf(t_2) & 0 & 0 \\ 0 & 0 & idf(t_3) & 0 \\ 0 & 0 & 0 & idf(t_4) \end{bmatrix}$$

$$= \begin{bmatrix} tf(t_1, d_1) \times idf(t_1) & tf(t_2, d_1) \times idf(t_2) & tf(t_3, d_1) \times idf(t_3) & tf(t_4, d_1) \times idf(t_4) \\ tf(t_1, d_2) \times idf(t_1) & tf(t_2, d_2) \times idf(t_2) & tf(t_3, d_2) \times idf(t_3) & tf(t_4, d_2) \times idf(t_4) \end{bmatrix}$$

Let's see now a concrete example of this multiplication:

$$M_{tf-idf} = M_{train} \times M_{idf} =$$

$$\begin{bmatrix} 0 & 1 & 1 & 1 \\ 0 & 2 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 0.69314718 & 0 & 0 & 0 \\ 0 & -0.40546511 & 0 & 0 \\ 0 & 0 & -0.40546511 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & -0.40546511 & -0.40546511 & 0 \\ 0 & -0.81093022 & -0.40546511 & 0 \end{bmatrix}$$

And finally, we can apply our L2 normalization process to the M_{tf-idf} matrix. Please note that this normalization is **"row-wise"** because we're going to handle each row of the matrix as a separated vector to be normalized, and not the matrix as a whole:

$$M_{tf-idf} = \frac{M_{tf-idf}}{\|M_{tf-idf}\|_2} = \begin{bmatrix} 0 & -0.70710678 & -0.70710678 & 0 \\ 0 & -0.89442719 & -0.4472136 & 0 \end{bmatrix}$$

And that is our pretty normalized tf-idf weight of our testing document set, which is actually a collection of unit vectors. If you take the L2-norm of each row of the matrix,

you'll see that they all have a L2-norm of 1.

Python practice

Environment Used: Python v.2.7.2, Numpy 1.6.1, Scipy v.0.9.0, Sklearn (Scikits.learn) v.0.9.

Now the section you were waiting for ! In this section I'll use Python to show each step of the tf-idf calculation using the **Scikit.learn** feature extraction module.

The first step is to create our training and testing document set and computing the term frequency matrix:

```
1. from sklearn.feature_extraction.text import CountVectorizer
2.
3. train_set = ("The sky is blue.", "The sun is bright.")
4. test_set = ("The sun in the sky is bright.",
5. "We can see the shining sun, the bright sun.")
6.
7. count_vectorizer = CountVectorizer()
8. count_vectorizer.fit_transform(train_set)
9. print "Vocabulary:", count_vectorizer.vocabulary
10.
11. # Vocabulary: {'blue': 0, 'sun': 1, 'bright': 2, 'sky': 3}
12.
13. freq_term_matrix = count_vectorizer.transform(test_set)
14. print freq_term_matrix.todense()
15.
16. #[[0 1 1 1]
17.  #[0 2 1 0]]
```

Now that we have the frequency term matrix (called **freq_term_matrix**), we can instantiate the **TfidfTransformer**, which is going to be responsible to calculate the tf-idf weights for our term frequency matrix:

```

1. from sklearn.feature_extraction.text import TfidfTransformer
2.
3. tfidf = TfidfTransformer(norm="l2")
4. tfidf.fit(freq_term_matrix)
5.
6. print "IDF:", tfidf.idf_
7.
8. # IDF: [ 0.69314718 -0.40546511 -0.40546511  0.      ]

```

Note that I've specified the norm as L2, this is optional (actually the default is L2-norm), but I've added the parameter to make it explicit to you that it's going to use the L2-norm. Also note that you can see the calculated idf weight by accessing the internal attribute called **idf_**. Now that **fit()** method has calculated the idf for the matrix, let's transform the **freq_term_matrix** to the tf-idf weight matrix:

```

1. tf_idf_matrix = tfidf.transform(freq_term_matrix)
2. print tf_idf_matrix.todense()
3.
4. # [[ 0.      -0.70710678 -0.70710678  0.      ]
5. # [ 0.      -0.89442719 -0.4472136  0.      ]]

```

And that is it, the **tf_idf_matrix** is actually our previous M_{tf-idf} matrix. You can accomplish the same effect by using the **Vectorizer** class of the Scikit.learn which is a vectorizer that automatically combines the **CountVectorizer** and the **TfidfTransformer** to you. See [this example](#) to know how to use it for the text classification process.

I really hope you liked the post, I tried to make it simple as possible even for people without the required mathematical background of linear algebra, etc. In the next Machine Learning post I'm expecting to show how you can use the tf-idf to calculate the cosine similarity.

If you liked it, feel free to comment and make suggestions, corrections, etc.

References

[Understanding Inverse Document Frequency: on theoretical arguments for IDF](#)

[Wikipedia :: tf-idf](#)

[The classic Vector Space Model](#)

[Sklearn text feature extraction code](#)

Updates

13 Mar 2015 – *Formating, fixed images issues.*

03 Oct 2011 – *Added the info about the environment used for Python examples*

← *Machine Learning :: Text feature extraction (tf-idf) – Part I*
C++11 user-defined literals and some constructions →

58 thoughts on “Machine Learning :: Text feature extraction (tf-idf) – Part II”



Severtcev

Wow!

Perfect intro in tf-idf, thank you very much! Very interesting, I've wanted to study this field for a long time and you posts it is a real gift. It would be very interesting to read more about use-cases of the technique. And may be you'll be interested, please, to shed some light on other methods of text corpus representation, if they exists? (sorry for bad English, I'm working to improve it, but there is still a lot of job to do)

↩ Reply



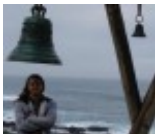
ogrisel

03/10/2011 at 22:34

Excellent work Christian! I am looking forward to reading your next posts on document classification, clustering and topics extraction with Naive Bayes, Stochastic Gradient Descent, Minibatch-k-Means and Non Negative Matrix factorization 😊

Also, the documentation of scikit-learn is really poor on the text feature extraction part (I am the main culprit...). Don't hesitate to join the mailing list if you want to give a hand and improve upon the current situation.

↩ Reply



03/10/2011 at 23:18

Great thanks Olivier. I really want to help sklearn, I just have to get



★ Christian
S. Perone

some more time to do that, you guys have done a great work, I'm really impressed by the amount of algorithms already implemented in the lib, keep the good work !

↩ Reply

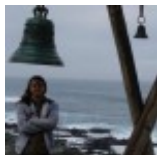


anand
jeyahar

04/10/2011 at 03:34

I like this tutorial better for the level of new concepts i am learning here.
That said, which version of scikits-learn are you using?.
The latest as installed by easy_install seems to have a different module hierarchy (i.e doesn't find feature_extraction in sklearn). If you could mention the version you used, i will just try out with those examples.

↩ Reply



★ Christian
S. Perone

04/10/2011 at 09:44

Hello Anand, I'm glad you liked it. I've added the information about the environment used just before the section "Python practice", I'm using the scikits.learn 0.9 (released a few weeks ago).

↩ Reply

Pingback: [Machine Learning :: Text feature extraction \(tf-idf\) – Part I | Pyevolve](#)

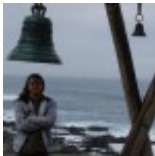


siamii

17/11/2011 at 01:57

Where's part 3? I've got to submit an assignment on Vector Space Modelling in 4 days. Any hope of putting it up over the weekend?

↩ Reply



★ Christian
S. Perone

18/11/2011 at 08:29

I've no date to publish it since I haven't got any time to write it =(

↩ Reply



Niu

24/11/2011 at 14:16

Thanks again for this complete and explicit tutorial and I am waiting for the coming section.

↩ Reply



Jason Wu

13/12/2011 at 00:05

Thanks Christian! a very nice work on vector space with sklearn. I just have one question, suppose I have computed the 'tf_idf_matrix', and I would like to compute the pair-wise cosine similarity (between each rows). I was having problem with the sparse matrix format, can you please give an example on that? Also my matrix is pretty big, say 25k by

60k. Thanks a lot!

↩ Reply



Khalid

10/01/2012 at 12:05

Great post... I understand what tf-idf and how to implement it with a concrete example. But I caught 2 things that I'm not sure about:

- 1- You called the 2 dimensional matrix M_{train} , but it has the tf values of the D3 and D4 documents, so you should've called that matrix M_{test} instead of M_{train} . Because D3 and D4 are our test documents.
- 2- When you calculate the idf value for the t_2 (which is 'sun') it should be $\log(2/4)$. Because number of the documents is 2. D3 has the word 'sun' 1 time, D4 has it 2 times. Which makes it 3 but we also add 1 to that value to get rid of divided by 0 problem. And this makes it 4... Am I right or am I missing something?

Thank you.

↩ Reply



arzu

28/01/2012 at 08:23

thanks... excellent post...

↩ Reply



04/02/2012 at 12:05



Jack

excellent post!

I have some question. From the last tf-idf weight matrix, how can we get the importance of term respectively(e.g. which is the most important term?). How can we use this matrix to classify documents

↩ Reply



Thanuj

13/02/2012 at 07:12

Thank You So Much. You explained it in such a simple way. It was really useful. Once again thanks a lot.

↩ Reply



Thanuj

13/02/2012 at 07:15

I have same doubt as Jack(last comment). From the last tf-idf weight matrix, how can we get the importance of term respectively(e.g. which is the most important term?). How can we use this matrix to classify documents.

↩ Reply



tintin

11/03/2012 at 07:12

I have a question..

After the tf-idf operation, we get a numpy array with values. Suppose we need to get the highest 50 values from the array. How can we do

that?

↩ Reply



Vikram
Bakhtiani

23/03/2012 at 12:12

Hey,

Thanx fr d code..was very helpful indeed !

1.For document clustering,after calculating inverted term frequency, shud i use any associativity coefficient like Jaccards coefficient and then apply the clustering algo like k-means or shud i apply d k-means directly to the document vectors after calculating inverted term frequency ?

2. How do u rate inverted term frequency for calcuating document vectors for document clustering ?

Thanks a ton fr the forth coming reply! 😊

↩ Reply



Frédérique
Passot

13/06/2012 at 16:23

@Khalid: what you're pointing out in 1- got me confused too for a minute (M_train vs M_test). I think you are mistaken on your second point, though, because what we are using is really the number of documents in which a term occurs, regardless of the number of times the term occurs in any given document. In this case, then, the

denominator in the idf value for t2 (“sun”) is indeed 2+1 (2 documents have the term “sun”, +1 to avoid a potential zero division error).

I’d love to read the third installment of this series too! I’d be particularly interested in learning more about feature selection. Is there an idiomatic way to get a sorted list of the terms with the highest tf.idf scores? How would you identify those terms overall? How would you get the terms which are the most responsible for a high or low cosine similarity (row by row)?

Thank you for the _great_ posts!

↩ Reply



Matthys
Meintjes

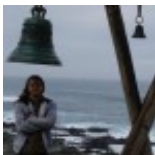
20/07/2012 at 17:20

Excellent article and a great introduction to td-idf normalization.

You have a very clear and structured way of explaining these difficult concepts.

Thanks!

↩ Reply



20/07/2012 at 22:52

Thanks for the feedback Matthys, I’m glad you liked the tutorial series.

★ Christian

S. Perone

↩ Reply

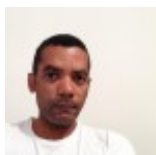


param

14/04/2015 at 09:39

very good & infomative tutorial.... please upload more tutorials related to documents clustering process.

↩ Reply



Laurent

27/07/2012 at 08:56

Excellent article ! Thank you Christian. You did a great job.

↩ Reply



Gavin Igor

24/08/2012 at 23:40

Can you provide any reference for doing cosine similarity using tfidf so we have the matrix of tf-idf how can we use that to calculate cosine.
Thanks for fantastic article.

↩ Reply



27/08/2012 at 17:48

Thanks so much for this and for explaining the whole tf-idf thing thoroughly.

lavender

↩ Reply

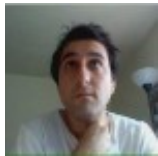


27/08/2012 at 18:35

Thanks for the feedback, I'm glad you liked the tutorial series.

★ Christian
S. Perone

↩ Reply



berkay

29/10/2012 at 09:24

Please correct me if i'm wrong
the formula after starting with "frequency we have calculated in the
first tutorial:" should Mtest not Mtrain. also after starting 'These idf
weights can be represented by a vector as:" should be idf_test not
idf_train.

Btw great series, can you give an simple approach for how to
implement classification?

↩ Reply



Divya

09/11/2012 at 02:15

Excellent it really helped me get through the concept of VSM and tf-idf.
Thanks Christian

↩ Reply



13/12/2012 at 12:57

Very good post. Congrats!!

Sergio

Showing your results, I have a question:

I read in the wikipedia:

The tf-idf value increases proportionally to the number of times a word appears in the document, but is offset by the frequency of the word in the corpus, which helps to control for the fact that some words are generally more common than others.

When I read it, I understand that if a word apperars in all documents is less important that a word that only appears in one document:

However, in the results, the word “sun” or “bright” are most important than “sky”.

I’m not sure of understand it completly.

↩ Reply

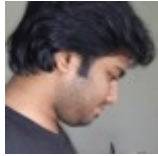


19/01/2013 at 17:00

Awesome! Explains TF-IDF very well. Waiting eagerly for your next post.

Bonson
Mampilli

↩ Reply

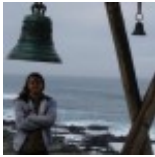


rahul
agarwal

04/07/2013 at 09:10

awesome work with a clear cut explanation . Even a layman can easily understand the subject..

↩ Reply



★ Christian
S. Perone

04/07/2013 at 14:17

Great thanks for the feedback Rahul !

↩ Reply



Jeremie

09/09/2013 at 18:41

Hello,

The explanation is awesome. I haven't seen a better one yet. I have trouble reproducing the results. It might be because of some update of sklearn.

Would it be possible for you to update the code?

It seem that the formula for computing the tf-idf vector has changed a little bit. Is a typo or another formula. Below is the link to the source code.

https://github.com/scikit-learn/scikit-learn/blob/master/sklearn/feature_extraction/text.py#L954

Many thanks

↩ Reply



11/09/2013 at 20:44

Terrific! I was familiar with tf-idf before but I found your scikits examples helpful as I'm trying to learn that package.

Susan

↩ Reply



11/09/2013 at 21:18

I'm glad you liked Susan, thanks for the feedback !

★ Christian
S. Perone

↩ Reply



26/10/2013 at 01:21

Thank you for writing such a detailed post. I learn allot.

Aziz

↩ Reply



13/11/2013 at 17:59

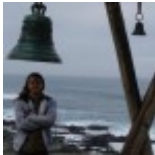
Excellent post! Stumbled on this by chance looking for more

Eugene
Chinveerap
han

information on CountVectorizer, but I'm glad I read through both of your posts (part 1 and part 2).

Bookmarking your blog now 😊

↩ Reply



★ Christian
S. Perone

14/11/2013 at 08:20

Great thanks for the feedback Eugene, I'm really glad you liked the tutorial series.

↩ Reply



me

25/02/2014 at 17:17

Does not seem to fit_transform() as you describe..
Any idea why?

```
>>> ts
```

```
('The sky is blue', 'The sun is bright')
```

```
>>> v7 = CountVectorizer()
```

```
>>> v7.fit_transform(ts)
```

```
<2x2 sparse matrix of type "
```

```
with 4 stored elements in COOrdinate format>
```

```
>>> print v7.vocabulary_
```

```
{u'is': 0, u'the': 1}
```

↩ Reply



Ash

20/06/2015 at 16:16

Actually, there are two small errors in the first Python sample.

1. CountVectorizer should be instantiated like so:

```
1 count_vectorizer = CountVectorizer(stop_words='english')
```

This will make sure the 'is', 'the' etc are removed.

2. To print the vocabulary, you have to add an underscore at the end.

```
1 print "Vocabulary:", count_vectorizer.vocabulary_
```

Excellent tutorial, just small things. hope it helps others.

↩ Reply



Kelvin John

28/04/2014 at 03:56

I loved your article.

↩ Reply



John

15/07/2014 at 15:54

I'm using scikit learn v.14. Is there any reason my results for running the exact same code would result in different results?

Anderton

↩ Reply



17/07/2014 at 13:04

Thanks for taking time to write this article. Found it very useful.

Karthik

↩ Reply



20/08/2014 at 07:41

Its useful.....thank you explaining the TD_IDF very elaborately..

Vijay

↩ Reply



11/09/2014 at 08:57

Thanks for the great explanation.

Mike

I have a question about calculation of the $\text{idf}(t\#)$.

In the first case, you wrote $\text{idf}(t_1) = \log(2/1)$, because we don't have such term in our collection, thus, we add 1 to the denominator. Now, in case t_2 , you wrote $\log(2/3)$, why the denominator is equal to 3 and not to 4 ($=1+2+1$)? In case t_3 , you write: $\log(2/3)$, thus the denominator is equal 3 ($=1+1+1$). I see here kind of inconsistency. Could you, please, explain, how did calculate the denominator value.

Thanks.

↩ Reply



★ Christian
S. Perone

11/09/2014 at 16:37

Hello Mike, thanks for the feedback. You're right, I just haven't fixed it yet due to lack of time to review it and recalculate the values.

↩ Reply



xpsycho

06/05/2015 at 15:58

You got it wrong, in the denominator you don't put the sum of the term in each document, you just sum all the documents that have at least one aparition of the term.

↩ Reply



huda

01/11/2014 at 03:58

This is good post

↩ Reply



huda

01/11/2014 at 04:16

it is good if you can provide way to know how use ft-idf in classification of document. I see that example (python code) but if there is algorithm that is best because no all people can understand this language.

Thanks

↩ Reply



13/11/2014 at 14:17

Great post, really helped me understand the tf-idf concept!

Ganesh

↩ Reply



24/12/2014 at 20:57

Nice Post

Samuel
Kahn

↩ Reply



22/01/2015 at 19:06

Nice. An explanation helps put things into perspective. Is tf-idf a good way to do clustering (e.g. use Jaccard analysis or variance against the average set from a known corpus)?

Andrew
Evans

Keep writing:)

↩ Reply



Neethu
Prem

25/01/2015 at 15:23

Hi Christian,

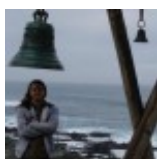
It makes me very excited and lucky to have read this article. The clarity of your understanding reflects in the clarity of the document. It makes me regain my confidence in the field of machine learning.

Thanks a ton for the beautiful explanation.

Would like to read more from you.

Thanks,
Neethu

↩ Reply



★ Christian
S. Perone

25/01/2015 at 23:49

Great thanks for the kind words Neethu ! I'm very glad you liked the tutorial series.

↩ Reply



esra'a ok

09/03/2015 at 18:01

thank you very very much, very wonderful and useful.

↩ Reply



09/03/2015 at 21:23

Thanks for the feedback Esra'a.

★ Christian
S. Perone

↩ Reply

Pingback: [Machine Learning :: Cosine Similarity for Vector Space Models \(Part III\) | Pyevolve](#)



29/03/2015 at 07:53

Thank you for the good wrap up. You mention a number of papers which compare L1 and L2 norm, I plan to study that a bit more in depth. You still know their names?

Arne

↩ Reply



27/04/2015 at 07:54

how can i calculate tf idf for my own text file which is located somewhere in my pc?

seher

↩ Reply

Leave a Reply

Your email address will not be published.

Name

Email

Website

Comment

You may use these HTML tags and attributes: ` <abbr title=""> <acronym title=""> <blockquote cite=""> <cite> <code> <del datetime=""> <i> <q cite=""> <s> <strike> `

Post Comment

Search ...

Meta

[Log in](#)

Tags

3d announce aqua arduino **Article**
benford's law bitcoin C c++0x
comparative evolution
evolutionary algorithms ga
genetic algorithm **Genetic**
Algorithms genetic
programming jit jython
machine learning matplotlib
modis **News** Philosophy
programming **Pyevolve**
Python raspberry pi rpi
Science sigevolution sklearn
sony psp stallion statistics svn terra
text mining tf-idf tsp twitter ubigraph
usgs vector space model visualization
vsm

[Entries RSS](#)

[Comments RSS](#)

[WordPress.org](#)

Proudly powered by [WordPress](#)