

Assignment

What does tf-idf mean?

Tf-idf stands for *term frequency-inverse document frequency*, and the tf-idf weight is a weight often used in information retrieval and text mining. This weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus. The importance 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. Variations of the tf-idf weighting scheme are often used by search engines as a central tool in scoring and ranking a document's relevance given a user query.

One of the simplest ranking functions is computed by summing the tf-idf for each query term; many more sophisticated ranking functions are variants of this simple model.

Tf-idf can be successfully used for stop-words filtering in various subject fields including text summarization and classification.

How to Compute:

Typically, the tf-idf weight is composed by two terms: the first computes the normalized Term Frequency (TF), aka. the number of times a word appears in a document, divided by the total number of words in that document; the second term is the Inverse Document Frequency (IDF), computed as the logarithm of the number of the documents in the corpus divided by the number of documents where the specific term appears.

- **TF:** Term Frequency, which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length (aka. the total number of terms in the document) as a way of normalization:

$$TF(t) = \frac{\text{Number of times term } t \text{ appears in a document}}{\text{Total number of terms in the document}}.$$

- **IDF:** Inverse Document Frequency, which measures how important a term is. While computing TF, all terms are considered equally important. However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:

$$IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents with term } t \text{ in it}} \cdot \text{for numerical stability we will be changing this formula little bit}$$

$$IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents with term } t \text{ in it} + 1} \cdot$$

Example

Consider a document containing 100 words wherein the word cat appears 3 times. The term frequency (i.e., tf) for cat is then $(3 / 100) = 0.03$. Now, assume we have 10 million documents and the word cat appears in one thousand of these. Then, the inverse document frequency (i.e., idf) is calculated as $\log(10,000,000 / 1,000) = 4$. Thus, the Tf-idf weight is the product of these quantities: $0.03 * 4 = 0.12$.

Task-1

1. Build a TFIDF Vectorizer & compare its results with Sklearn:

- As a part of this task you will be implementing TFIDF vectorizer on a collection of text documents.
- You should compare the results of your own implementation of TFIDF vectorizer with that of sklearn's implementation TFIDF vectorizer.
- Sklearn does few more tweaks in the implementation of its version of TFIDF vectorizer, so to replicate the exact results you would need to add following things to your custom implementation of tfidf vectorizer:
 1. Sklearn has its vocabulary generated from idf sorted in alphabetical order
 2. Sklearn formula of idf is different from the standard textbook formula. Here the constant "1" is added to the numerator and denominator of the idf as if an extra document was seen containing every term in the collection exactly once, which prevents zero divisions.

$$IDF(t) = 1 + \log_e \frac{1 + \text{Total number of documents in collection}}{1 + \text{Number of documents with term } t \text{ in it}} \cdot$$
 3. Sklearn applies L2-normalization on its output matrix.
 4. The final output of sklearn tfidf vectorizer is a sparse matrix.
- Steps to approach this task:
 1. You would have to write both fit and transform methods for your custom implementation of tfidf vectorizer.

2. Print out the alphabetically sorted voacb after you fit your data and check if its the same as that of the feature names from sklearn tfidf vectorizer.
3. Print out the idf values from your implementation and check if its the same as that of sklearn's tfidf vectorizer idf values.
4. Once you get your voacb and idf values to be same as that of sklearn's implementation of tfidf vectorizer, proceed to the below steps.
5. Make sure the output of your implementation is a sparse matrix. Before generating the final output, you need to normalize your sparse matrix using L2 normalization. You can refer to this link <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html>
6. After completing the above steps, print the output of your custom implementation and compare it with sklearn's implementation of tfidf vectorizer.
7. To check the output of a single document in your collection of documents, you can convert the sparse matrix related only to that document into dense matrix and print it.

Note-1: All the necessary outputs of sklearn's tfidf vectorizer have been provided as reference in this notebook, you can compare your outputs as mentioned in the above steps, with these outputs.

Note-2: The output of your custom implementation and that of sklearn's implementation would match only with the collection of document strings provided to you as reference in this notebook. It would not match for strings that contain capital letters or punctuations, etc, because sklearn version of tfidf vectorizer deals with such strings in a different way. To know further details about how sklearn tfidf vectorizer works with such string, you can always refer to its official documentation.

Note-3: During this task, it would be helpful for you to debug the code you write with print statements wherever necessary. But when you are finally submitting the assignment, make sure your code is readable and try not to print things which are not part of this task.

Corpus

```
In [1]: ## SkLearn# Collection of string documents

corpus = [
    'this is the first document',
    'this document is the second document',
    'and this is the third one',
    'is this the first document',
]
```

SkLearn Implementation

```
In [2]: from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer()
vectorizer.fit(corpus)
skl_output = vectorizer.transform(corpus)
```

```
In [3]: # sklearn feature names, they are sorted in alphabetic order by default.
```

```
print(vectorizer.get_feature_names())
```

```
['and', 'document', 'first', 'is', 'one', 'second', 'the', 'third', 'this']
```

```
In [ ]: # Here we will print the sklearn tfidf vectorizer idf values after applying the fit method
# After using the fit function on the corpus the vocab has 9 words in it, and each has its idf value.
```

```
print(vectorizer.idf_)
```

```
[1.91629073 1.22314355 1.51082562 1.          1.91629073 1.91629073
 1.          1.91629073 1.          ]
```

```
In [ ]: # shape of sklearn tfidf vectorizer output after applying transform method.
```

```
skl_output.shape
```

```
Out[ ]: (4, 9)
```

```
In [ ]: # sklearn tfidf values for first line of the above corpus.
# Here the output is a sparse matrix
```

```
print(skl_output[0])
```

```
(0, 8)      0.38408524091481483
(0, 6)      0.38408524091481483
(0, 3)      0.38408524091481483
(0, 2)      0.5802858236844359
(0, 1)      0.46979138557992045
```

```
In [ ]: # sklearn tfidf values for first line of the above corpus.
# To understand the output better, here we are converting the sparse output matrix to dense matrix and printing it.
# Notice that this output is normalized using L2 normalization. sklearn does this by default.
```

```
print(skl_output[0].toarray())
```

```
[[0.          0.46979139 0.58028582 0.38408524 0.          0.  
  0.38408524 0.          0.38408524]]
```

Your custom implementation

```
In [168... # Write your code here.  
# Make sure its well documented and readable with appropriate comments.  
# Compare your results with the above sklearn tfidf vectorizer  
# You are not supposed to use any other library apart from the ones given below  
  
from collections import Counter  
from tqdm import tqdm  
from scipy.sparse import csr_matrix  
import math  
import operator  
from sklearn.preprocessing import normalize  
import numpy  
  
corpus = [  
    'this is the first document',  
    'this document is the second document',  
    'and this is the third one',  
    'is this the first document',  
]  
  
def fit(dataset): #returns dictionary of vocab with indexes , vocab with IDF values  
    unique_words = set() #set used to remove duplicacy  
    vocab = {} #IDF with Vocab  
    if isinstance (dataset , (list,)):  
        lst = []  
        lst1 = []  
        for row in dataset:  
            lst.append(set(row.split(" "))) #remove duplicates in a sentence  
  
            for i in lst:  
                lst1.append(" ".join(i)) #combine sets into list
```

```

lst1 = sorted(lst1) #sort list
new_dataset_idf = " ".join(lst1) # joining entire dataset

#To get the count of documents a word is present in
total_word_count_idf = dict(Counter(new_dataset_idf.split(" "))) # dict of boolean word count from each document

for row in dataset:    #for each row
    for word in row.split(" "):    #for each word in a row - split words on space
        unique_words.add(word)    #add each word to the set , duplicates will not be added in a set hence an
unique_words = sorted(list(unique_words))    #convert to list then sort

for word in unique_words:
    vocab[word] = (1+math.log((1+len(dataset))/(1+total_word_count_idf.get(word)))) #vocab of IDF
vocab_idx = {word:idx for idx , word in enumerate(unique_words)} #Vocab with index

return vocab , vocab_idx

def transform(dataset,vocab):
    rows = []
    columns = []
    values = []
    res = fit(dataset)
    idf = res[0] #vocab with idf
    vocab_idx = res[1] #vocab with index

    if isinstance (dataset , (list,)):

        for idx , row in enumerate(dataset):

            word_freq = dict(Counter(row.split(" "))) #dictionary of a word as key and frequency

            temp_set = set()
            for word in row:
                temp_set.add(word)
            Total_words = len(temp_set) #distinct count of words in every text item
            for word, freq in word_freq.items(): # for every item in the dictionary

```

```

        col_index = vocab_idx.get(word , -1) #ignore the words not present in vocab

        if col_index != -1:
            rows.append(idx)          #row index
            columns.append(col_index) #column index
            values.append((freq/Total_words)*float(idf.get(word))) #final TFIDF value

    return csr_matrix((values, (rows,columns)),shape = (len(dataset),len(vocab))) #returns sparse matrix

vocab = fit(corpus)[0] #vocab with indexes
output = normalize(transform(corpus, vocab)) #normalized output of transform to match sklearn's implementation - l2
print("Sparse Matrix-----")
print(" ")
print(output[0]) #sparse output for the first line item
print(" ")
print("Array-----")
print(" ")
print(output[0].toarray()) #Dense output for the first line item

```

Sparse Matrix-----

```

(0, 1)      0.46979138557992045
(0, 2)      0.5802858236844359
(0, 3)      0.3840852409148149
(0, 6)      0.3840852409148149
(0, 8)      0.3840852409148149

```

Array-----

```

[[0.         0.46979139 0.58028582 0.38408524 0.         0.
  0.38408524 0.         0.38408524]]

```

Task-2

2. Implement max features functionality:

- As a part of this task you have to modify your fit and transform functions so that your vocab will contain only 50 terms with top idf scores.

- This task is similar to your previous task, just that here your vocabulary is limited to only top 50 features names based on their idf values. Basically your output will have exactly 50 columns and the number of rows will depend on the number of documents you have in your corpus.
- Here you will be given a pickle file, with file name **cleaned_strings**. You would have to load the corpus from this file and use it as input to your tfidf vectorizer.
- Steps to approach this task:
 1. You would have to write both fit and transform methods for your custom implementation of tfidf vectorizer, just like in the previous task. Additionally, here you have to limit the number of features generated to 50 as described above.
 2. Now sort your vocab based in descending order of idf values and print out the words in the sorted vocab after you fit your data. Here you should be getting only 50 terms in your vocab. And make sure to print idf values for each term in your vocab.
 3. Make sure the output of your implementation is a sparse matrix. Before generating the final output, you need to normalize your sparse matrix using L2 normalization. You can refer to this link <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html>
 4. Now check the output of a single document in your collection of documents, you can convert the sparse matrix related only to that document into dense matrix and print it. And this dense matrix should contain 1 row and 50 columns.

```
In [169... # Below is the code to load the cleaned_strings pickle file provided
# Here corpus is of list type

import pickle
with open('cleaned_strings', 'rb') as f:
    corpus = pickle.load(f)

# printing the length of the corpus loaded
print("Number of documents in corpus = ", len(corpus))
```

Number of documents in corpus = 746

```
In [170... # Write your code here.
# Make sure its well documented and readable with appropriate comments.
# Compare your results with the above sklearn tfidf vectorizer
# You are not supposed to use any other library apart from the ones given below
```



```

from collections import Counter
from tqdm import tqdm
from scipy.sparse import csr_matrix
import math
import operator
from sklearn.preprocessing import normalize
import numpy

def fit(dataset): #returns dictionary of vocab with indexes , vocab with IDF values
    unique_words = set() #set used to remove duplicacy
    vocab = {} #IDF with Vocab
    new_vocab = {}
    if isinstance (dataset , (list,)):
        lst = []
        lst1 = []
        for row in dataset:
            lst.append(set(row.split(" "))) #remove duplicates in a sentence

        for i in lst:
            lst1.append(" ".join(i)) #combine sets into list

        lst1 = sorted(lst1) #sort list
        new_dataset_idf = " ".join(lst1) # joining entire dataset

        #To get the count of documents a word is present in
        total_word_count_idf = dict(Counter(new_dataset_idf.split(" "))) # dict of boolean word count from each document

        for row in dataset: #for each row
            for word in row.split(" "): #for each word in a row - split words on space
                unique_words.add(word) #add each word to the set , duplicates will not be added in a set hence an
            unique_words = sorted(list(unique_words)) #convert to list then sort

        for word in unique_words:
            vocab[word] = (1+math.log((1+len(dataset))/(1+total_word_count_idf.get(word)))) #vocab of IDF
        new_vocab = {k: v for k, v in sorted(vocab.items(), reverse=True, key=operator.itemgetter(1))[:50]} #select top 50 words

        new_unique_words = []
        new_unique_words = new_vocab.keys()

```

```

new_unique_words = sorted(new_unique_words)
vocab_idx = {word:idx for idx, word in enumerate(new_unique_words)} #Vocab with index

return new_vocab, vocab_idx

def transform(dataset,vocab):
    rows = []
    columns = []
    values = []
    res = fit(dataset)
    idf = res[0] #vocab with idf
    vocab_idx = res[1] #vocab with index

    print("Top 50 IDF Values")
    print(idf)
    if isinstance (dataset, (list,)):

        for idx, row in enumerate(dataset):

            word_freq = dict(Counter(row.split(" "))) #dictionary of a word as key and frequency

            temp_set = set()
            for word in row:
                temp_set.add(word)
            Total_words = len(temp_set) #distinct count of words in every text item
            for word, freq in word_freq.items(): # for every item in the dictionary

                col_index = vocab_idx.get(word, -1) #ignore the words not present in vocab

                if col_index != -1:
                    rows.append(idx) #row index
                    columns.append(col_index) #column index
                    values.append((freq/Total_words)*float(idf.get(word))) #final TFIDF value

    return csr_matrix((values, (rows,columns)),shape = (len(dataset),len(vocab))) #returns sparse matrix

vocab = fit(corpus)[0] #vocab with indexes
output = normalize(transform(corpus, vocab)) #normalized output of transform to match sklearn's implementation - l2
print("Sparse Matrix-----")

```

```
print(" ")
print(output[0]) #sparse output for the first line item
print(" ")
print("Array-----")
print(" ")
print(output[0].toarray()) #Dense output for the first line item (1, 50)
```

Top 50 IDF Values

'aailiyah': 6.922918004572872, 'abandoned': 6.922918004572872, 'abroad': 6.922918004572872, 'abstruse': 6.922918004572872, 'academy': 6.922918004572872, 'accents': 6.922918004572872, 'accessible': 6.922918004572872, 'acclaimed': 6.922918004572872, 'accolades': 6.922918004572872, 'accurate': 6.922918004572872, 'accurately': 6.922918004572872, 'achille': 6.922918004572872, 'ackerman': 6.922918004572872, 'actions': 6.922918004572872, 'adams': 6.922918004572872, 'add': 6.922918004572872, 'added': 6.922918004572872, 'admins': 6.922918004572872, 'admiration': 6.922918004572872, 'admitted': 6.922918004572872, 'adrift': 6.922918004572872, 'adventure': 6.922918004572872, 'aesthetically': 6.922918004572872, 'affected': 6.922918004572872, 'affleck': 6.922918004572872, 'afternoon': 6.922918004572872, 'aged': 6.922918004572872, 'ages': 6.922918004572872, 'agree': 6.922918004572872, 'agreed': 6.922918004572872, 'aimless': 6.922918004572872, 'aired': 6.922918004572872, 'akasha': 6.922918004572872, 'akin': 6.922918004572872, 'alert': 6.922918004572872, 'alike': 6.922918004572872, 'allison': 6.922918004572872, 'allow': 6.922918004572872, 'allowing': 6.922918004572872, 'alongside': 6.922918004572872, 'amateurish': 6.922918004572872, 'amaze': 6.922918004572872, 'amazed': 6.922918004572872, 'amazingly': 6.922918004572872, 'amusing': 6.922918004572872, 'amust': 6.922918004572872, 'anatomist': 6.922918004572872, 'angel': 6.922918004572872, 'angela': 6.922918004572872, 'angelina': 6.922918004572872}

Sparse Matrix-----

$(0, 30)$	1.0
-----------	-----

Array-----

[illegible]