

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:

$$IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents with term } t \text{ in it}}. \text{ for numerical stability we will be changing this formula little bit } IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents with term } t \text{ in it} + 1}.$$

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 of 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. $IDF(t) = 1 + \log_e \frac{1 + \text{Total number of documents in collection}}{1 + \text{Number of documents with term } t \text{ in it}}$.
 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 vocab 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 vocab 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 [4]: # Here we will print the sklearn tfidf vectorizer idf values after applying t.
# After using the fit function on the corpus the vocab has 9 words in it, and

print(vectorizer.idf_)

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

```
In [5]: # shape of sklearn tfidf vectorizer output after applying transform method.

skl_output.shape
```

Out[5]: (4, 9)

```
In [6]: # 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 [7]: # sklearn tfidf values for first line of the above corpus.
# To understand the output better, here we are converting the sparse output m
# Notice that this output is normalized using L2 normalization. sklearn does

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

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

Your custom implementation

```
In [8]: # 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
import re

# ---- CONSTANTS ----

VOCABULARY = 'Vocabulary'
FEATURES = 'Features'
IDFS = 'InverseDocumentFrequency'

# ---- UTILS ----

def sent_to_words(sentence):
    """
    Convert sentence to List of Words

    #!! Currently using simple space to split words but can be enhanced with
    #TODO clean more for better nuance
    pattern = r"\s"
    w1 = re.split(pattern, sentence)
    w2 = [w for w in w1 if len(w1) > 2]
    return w2

def calc_idf(vocab, dataset, *, prec=2):
    """
    calculate the idf ie Inverse Document Frequency from given vocab & dataset

    :param: vocab :- vocabulary ie word to column
    :param: dataset :- list of documents
    :param: prec :- precision for floating val of idf

    :return: float :- idf value mapping for each word in vocabulary
    """
    doc_cnts = len(dataset)
```

```

log_doc_cnts = math.log1p(doc_cnts)
idf = {}

# creating the list of words for each doc
corpus = [sent_to_words(d) for d in dataset]

for word in vocab:
    # number of documents with term t in it
    doc_freq = sum([word in doc for doc in corpus])
    inverse_doc_freq = 1 + log_doc_cnts - math.log1p(doc_freq) #IDF(word)
    idf[word] = round(inverse_doc_freq, prec)
return idf

# --- IMPL ----

def fit(dataset):
    """
    Generate Vocabulary from given Corpus (ie collection of documents)
    Vocabulary maps the word (ie feature) to its column number

    :return: dictionary of Vocabulary & Features
    """
    if isinstance(dataset, (list,)):
        unique_words = {w for s in dataset for w in sent_to_words(s)}
        words_list = sorted(unique_words)
        vocab = {j:i for i,j in enumerate(words_list)} # word -> column

        return {VOCABULARY: vocab, FEATURES: words_list}
    else:
        print("you need to pass list of sentence")

def transform(dataset, vocab, idfs):
    """
    Transform each documents in list to a new Sparse Vector repr for given vo

    :param dataset: - list of documents
    :param vocab: - mapping of words to its counn number in sparse matrix repr
    :param idfs: - mapping of words to its idf values

    :return: list of sparse vectors corresp to each document
    """

    # (r, c, v) where r = row num; c = col num; v = tfidf vals
    rows, columns, vals = [], [], []

    if isinstance(dataset, (list,)): # check if dataset is `List of List`
        for row, doc in enumerate(tqdm(dataset)):
            # 1 Find the Sparse Features
            # ie sparse features for document {doc} := common words between d
            words = sent_to_words(doc)
            unique_words = set(words)
            common_words = vocab.keys() & unique_words # sparse features

            # 2 calculate the TF ie term frequency for sparse words in docume
            tf = {w: words.count(w) for w in common_words}

            # 3 Calculate TfIdf ie word -> tf*idf
            tfidf = {w: tf[w] * idfs.get(w, 0) for w in tf}

            # 4 Calculate L2 norm (tfidf.values() -> will hold all sparse va
            l2_norm = math.sqrt(sum([x**2 for x in tfidf.values()]))

            # 5 Prepare Sparse Vector Repr for {doc}
            for w in tf:

```

```

        rows.append(row)                # row number
        columns.append(vocab.get(w, -1)) # column number
        vals.append(tfidf[w]/l2_norm)    # Tf-Idf value

    return csr_matrix((vals, (rows,columns)), shape=(len(dataset),len(vocab)))

```

```

In [9]: # 1. Fit the Data
        d = fit(corpus)
        vocab = d[VOCABULARY] # word -> col_no

        # 2. Find the IDF
        idfs = calc_idf(vocab, corpus, prec=6)

        # 3. Find the Sparse Features for Data Points
        sparse_m = transform(corpus, vocab, idfs)

```

```
100%|██████████| 4/4 [00:00<00:00, 613.67it/s]
```

```

In [10]: print('Features : \n', d[FEATURES])
        print('Shape : ', sparse_m.shape)

Features :
['and', 'document', 'first', 'is', 'one', 'second', 'the', 'third', 'this']
Shape : (4, 9)

```

```

In [11]: sparse_m[0]

```

```

Out[11]: <1x9 sparse matrix of type '<class 'numpy.float64'>'
         with 5 stored elements in Compressed Sparse Row format>

```

```

In [12]: sparse_m[0].toarray()

```

```

Out[12]: array([[0.          , 0.46979148, 0.58028587, 0.38408518, 0.          ,
                  0.          , 0.38408518, 0.          , 0.38408518]])

```

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 voacb 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 [13]: # 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 [14]: # Write your code here.
# Try not to hardcode any values.
# Make sure its well documented and readable with appropriate comments.
```

```
In [15]: # --- UTIL ---
def get_top_features(n, vocab, idfs):
    """
    calculates top n features & corresponding vocab & idfs
    (features are considered top with increasing val of their idf values)

    : param n: - top n words to pick s.t.theirs idfs are high
    : param vocab: - vocabulary of words to their column numbers
    : param idfs: - mapping of words to their idf values

    :return : - mapping of Features, Vocabulary & IDFs
    """
    top_50_words = sorted(vocab, key=idfs.__getitem__, reverse=True)[:50]
    new_vocab = {w:i for i,w in enumerate(top_50_words)}
    new_idfs = {w:idfs[w] for w in top_50_words}

    return {FEATURES: top_50_words,VOCABULARY: new_vocab, IDFS: new_idfs}
```

```
In [16]: # 1. Fit the Data on Entire Data
d = fit(corpus)
vocab = d[VOCABULARY]

# 2. Find the IDF
idfs = calc_idf(vocab, corpus, prec=6)

# 3. Get Top 50 Features
```

