```
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.
           </font>
          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) = rac{	ext{Number of times term t appears in a document}}{	ext{T}}
                         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 rac{	ext{Total number of documents}}{	ext{Number of documents with term t in it}}. for numerical stabiltiy we will be changing this formula little bit IDF(t) = \log_e rac{	ext{Total number of documents}}{	ext{Number of documents with term t in it}}.
          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.
            </font>
          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 sklearns 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 sroted 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 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 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 sklearns tfidf vectorizer idf
                  values.
                 4. Once you get your voacb and idf values to be same as that of sklearns 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 sklearns
                  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 sklearns 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 sklearns 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 [57]: ## 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
          from sklearn.feature_extraction.text import TfidfVectorizer
           vectorizer = TfidfVectorizer()
           vectorizer.fit(corpus)
           skl_output = vectorizer.transform(corpus)
           print(skl_output)
             (0, 8)
                             0.38408524091481483
             (0, 6)
                             0.38408524091481483
             (0, 3)
                             0.38408524091481483
             (0, 2)
                             0.5802858236844359
             (0, 1)
                             0.46979138557992045
             (1, 8)
                             0.281088674033753
             (1, 6)
                             0.281088674033753
             (1, 5)
                             0.5386476208856763
             (1, 3)
                             0.281088674033753
             (1, 1)
             (2, 8)
                             0.267103787642168
             (2, 7)
                             0.511848512707169
             (2, 6)
                             0.267103787642168
             (2, 4)
                             0.511848512707169
             (2, 3)
                             0.267103787642168
             (2, 0)
                             0.511848512707169
             (3, 8)
                             0.38408524091481483
             (3, 6)
                             0.38408524091481483
             (3, 3)
                             0.38408524091481483
             (3, 2)
                             0.5802858236844359
             (3, 1)
                             0.46979138557992045
In [59]: # 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 [60]: # 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 i
           df value.
           print(vectorizer.idf_)
                                                                1.91629073 1.91629073
           [1.91629073 1.22314355 1.51082562 1.
                         1.91629073 1.
In [61]: # shape of sklearn tfidf vectorizer output after applying transform method.
          skl_output.shape
Out[61]: (4, 9)
In [62]: # 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 [63]: # 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 defaul
           print(skl_output[0].toarray())
                          0.46979139 0.58028582 0.38408524 0.
             0.38408524 0.
                                       0.38408524]]
          Your custom implementation
In [64]: # Write your code here.
           # Make sure its well documented and readble 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
           import tqdm
           from scipy.sparse import csr_matrix
           import math
           import operator
           from sklearn.preprocessing import normalize
           import numpy
In [65]: corpus = [
                 'this is the first document',
                 'this document is the second document',
                 'and this is the third one',
                 'is this the first document',
          Fit Method
In [66]: def fit(data):
               #Initializing an empty set called uniquewords because set cannot contain duplicate items
               uniquewords=set()
               if isinstance(data,(list)):
               #Looping through each row in the corpus
                    for row in data:
                         for word in row.split(" "):
                 #if the length of the word is greater than 2 then add it to the set
                             if len(word)>=2:
                                  uniquewords.add(word)
               #Sorting the words alphabetically since sklearn sorts the list alphabetically
               uniquewords=sorted(list(uniquewords))
               #Creating a dictionary with key as the uniqueword and value as the dimension number
               vocab={j:i for i,j in enumerate(uniquewords)}
               idf_values=IDF(data,uniquewords)
               return vocab,idf_values
           def IDF(corpus,uniquewords): #Defining a separate function to calculate IDF values
               #Creating a dictionary for to store the values for each uniqueword
               idf_dict={}
               n=len(corpus)
               # Counting the frequency of each uniqueword in each row
               for i in uniquewords:
                    count=0
                    for row in corpus:
                         if i in row.split():
                              count=count+1
                    \#IDF(t)=1+loge((1 + Total number of documents in collection)/(1+Number of documents)
            with term t in it))
                    idf_dict[i]=1+(math.log((1+n)/(1+count)))
               return idf_dict
           vocab,idf_values=fit(corpus)
          print('Feature names-\n', vocab, "\n")
          Feature names-
           {'and': 0, 'document': 1, 'first': 2, 'is': 3, 'one': 4, 'second': 5, 'the': 6, 'third': 7,
           'this': 8}
In [67]: print('IDF-\n',idf_values)
           { 'and ': 1.916290731874155, 'document ': 1.2231435513142097, 'first ': 1.5108256237659907, 'i
          s': 1.0, 'one': 1.916290731874155, 'second': 1.916290731874155, 'the': 1.0, 'third': 1.916290
          731874155, 'this': 1.0}
          The keys and values in the dictionary are the unique words and idf values of the corpus and they are similar to the output given from
          sklearn implementation
          Transform Method
In [71]: from collections import Counter
           import tqdm
           def transform(data, vocab):
               #For Creating a sparse matrix, we need to store the index, dimension number and the tf_idf
           values of the document.
               rows=[]
               columns=[]
               values=[]
               #Checking if the corpus is a list
               if isinstance(data,(list)):
               #For each document in the corpus, we are creating a dictionary where word is the key and
           frequency of the word is the value
                    for index,row in enumerate((data)):
                         w_freq=dict(Counter(row.split()))
               #Each uniqueword in the corpus
                         for word, freq in w_freq.items():
                              if len(word)<2:</pre>
                                  continue
               #Checking if the uniqueword is in our vocabulary.If the word is present then return the
            value else return -1
                              column_index=vocab.get(word, -1)
                              if column_index!=-1:
                              #Store the index of the document
                                   rows.append(index)
                              #Store the dimension number of the document
                                   columns.append(column_index)
                               #Multiplying term frequency and inverse document frequency to get tf_idf va
           1ue
                                   tf_idf_value=((freq/len(row.split()))*idf_values[word])
                               #storing the tf_idf value in the list
                                  values.append(tf_idf_value)
                               #Creating a sparse matrix
                                  sparse_matrix=csr_matrix((values, (rows,columns)), shape=(len(data),len(
           vocab)))
                               #Applying the normalisation formula for the output
                                  normalized_out=normalize(sparse_matrix,norm='12',axis=1, copy=True, retu
           rn_norm=False)
                    return normalized_out
               else:
                    print("you need to pass list of strings")
           print('Normalised TF_IDF values:')
           print(transform(corpus, vocab), '\n')
           print('Shape of sparse matrix', transform(corpus, vocab).shape, '\n')
          Normalised TF_IDF values:
             (0, 1)
                             0.4697913855799205
             (0, 2)
                             0.580285823684436
             (0, 3)
                             0.3840852409148149
             (0, 6)
                             0.3840852409148149
             (0, 8)
                             0.3840852409148149
             (1, 1)
                             0.6876235979836937
             (1, 3)
                             0.2810886740337529
             (1, 5)
                             0.5386476208856762
             (1, 6)
                             0.2810886740337529
             (1, 8)
                             0.2810886740337529
             (2, 0)
                             0.511848512707169
             (2, 3)
                             0.267103787642168
             (2, 4)
                             0.511848512707169
                             0.267103787642168
             (2, 6)
             (2, 7)
                             0.511848512707169
             (2, 8)
                             0.267103787642168
             (3, 1)
                             0.4697913855799205
                              0.580285823684436
             (3, 2)
                             0.3840852409148149
             (3, 3)
             (3, 6)
                             0.3840852409148149
                             0.3840852409148149
             (3, 8)
          Shape of sparse matrix (4, 9)
In [74]: | print("Dense array-")
          print(transform(corpus, vocab)[0].toarray())
          Dense array-
          [[0.
                          0.46979139 0.58028582 0.38408524 0.
                                                                              0.
             0.38408524 0.
                                      0.38408524]]
          Fit and Transform methods have been performed on the corpus and the result obtained is similar to the result obtained from sklearn
          implementation
          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

             • 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 give 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 [75]: # 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:
               corpus1 = pickle.load(f)
           # printing the length of the corpus loaded
           print("Number of documents in corpus = ",len(corpus1))
          Number of documents in corpus = 746
          Fit Method
In [76]: import numpy as np
           def fit(data):
               #Initializing an empty set called uniquewords because set cannot contain duplicate items
               uniquewords=set()
               Idf=[]
               #Looping through each row in the corpus
               for row in data:
                    for word in row.split(" "):
                 #if the length of the word is greater than 2 then add it to the set
                         if len(word)>=2:
                             uniquewords.add(word)
               #Sorting the words alphabetically since sklearn sorts the list alphabetically
               uniquewords=sorted(list(uniquewords))
               #Creating a dictionary with key as the uniqueword and value as the dimension number
               idf_values=IDF(data,uniquewords)
               vocab=(dict(zip(uniquewords,idf_values)))
               vocab2={j:i for i,j in enumerate(vocab1)}
               return vocab2,idf_values
           def IDF(corpus,uniquewords): #Defining a separate function to calculate IDF values
               #Creating a dictionary for to store the values for each uniqueword
               idf_dict={}
               idf50={}
               top50_idf=[]
               n=len(corpus)
               # Counting the frequency of each uniqueword
               for i in uniquewords:
                    count=0
                    for row in corpus:
                         if i in row.split():
                              count=count+1
                         idf_dict[i]=1+(math.log((1+n)/(1+count)))
               #Sorting the list to get top50 idf values
               top50_idf=sorted(idf_dict.items(), key=lambda x:x[1], reverse=True)[:50]
               idf50=dict(top50_idf)
               return idf50
           vocab2,idf_values=fit(corpus1)
           print('Words containing top50 IDF values are:\n')
          print(idf_values, "\n")
          Words containing top50 IDF values are:
           {'aailiyah': 6.922918004572872, 'abandoned': 6.922918004572872, 'abroad': 6.922918004572872,
           'abstruse': 6.922918004572872, 'academy': 6.922918004572872, 'accents': 6.922918004572872, 'a
          ccessible': 6.922918004572872, 'acclaimed': 6.922918004572872, 'accolades': 6.92291800457287
          2, 'accurate': 6.922918004572872, 'accurately': 6.922918004572872, 'achille': 6.9229180045728
          72, 'ackerman': 6.922918004572872, 'actions': 6.922918004572872, 'adams': 6.922918004572872,
           'add': 6.922918004572872, 'added': 6.922918004572872, 'admins': 6.922918004572872, 'admiratio
          n': 6.922918004572872, 'admitted': 6.922918004572872, 'adrift': 6.922918004572872, 'adventur
          e': 6.922918004572872, 'aesthetically': 6.922918004572872, 'affected': 6.922918004572872, 'af
          fleck': 6.922918004572872, 'afternoon': 6.922918004572872, 'aged': 6.922918004572872, 'ages':
          6.922918004572872, 'agree': 6.922918004572872, 'agreed': 6.922918004572872, 'aimless': 6.9229
          18004572872, 'aired': 6.922918004572872, 'akasha': 6.922918004572872, 'akin': 6.9229180045728
          72, 'alert': 6.922918004572872, 'alike': 6.922918004572872, 'allison': 6.922918004572872, 'al
          low': 6.922918004572872, 'allowing': 6.922918004572872, 'alongside': 6.922918004572872, 'amat
          eurish': 6.922918004572872, 'amaze': 6.922918004572872, 'amazed': 6.922918004572872, 'amazing
          ly': 6.922918004572872, 'amusing': 6.922918004572872, 'amust': 6.922918004572872, 'anatomis
          t': 6.922918004572872, 'angel': 6.922918004572872, 'angela': 6.922918004572872, 'angelina':
          6.922918004572872}
          Transform Method
In [77]: from collections import Counter
           import tqdm
           #For Creating a sparse matrix, we need to store the index, dimension number and the tf_idf va
           lues of the document.
           def transform(data, vocab):
               rows=[]
               columns=[]
               values=[]
              #Checking if the corpus is a list
               if isinstance(data,(list)):
               #For each document in the corpus, we are creating a dictionary where word is the key and
           frequency of the word is the value
                    for index,row in enumerate((data)):
                         w_freq=dict(Counter(row.split()))
               #Each uniqueword in the corpus
                         for word, freq in w_freq.items():
                             if len(word)<2:</pre>
                                  continue
                          #Checking if the uniqueword is in our vocabulary. If the word is present then re
           turn the value else return -1
                              column_index=vocab.get(word, -1)
                              if column index!=-1:
                              #Store the index of the document
                                  rows.append(index)
                              #Store the dimension number of the document
                                  columns.append(column_index)
                               #Multiplying term frequency and inverse document frequency to get tf_idf va
           1ue
                                  tf_idf_value=((freq/len(row.split()))*idf_values[word])
                               #storing the tf_idf value in the list
                                  values.append(tf_idf_value)
                               #Creating a sparse matrix
                                  sparse_matrix=csr_matrix((values, (rows,columns)), shape=(len(data),len(
           vocab)))
                               #Applying the normalisation formula for the output
                                  normalized_out=normalize(sparse_matrix,norm='l2',axis=1, copy=True, retu
           rn_norm=False)
                    return normalized_out
```

else:

print(output)

(746, 50)

(0, 30)

(68, 24)

(72, 29)

(74, 31)

(119, 33)

(135, 3)

(135, 10)

(135, 18)
(135, 20)

(135, 36)

(135, 40)

(135, 41)
(176, 49)

(181, 13)

(192, 21)

(193, 23)

(216, 2)

(222, 47)

(225, 19)

(227, 17) (241, 44)

(270, 1)

(290, 25)

(333, 26)

(334, 15)

(341, 43)

(344, 42)

(348, 8)

(377, 37)

(409, 5)

(430, 39)

(457, 45)

(461, 4)

(465, 38)

(475, 35)

(493, 6) (500, 48) (548, 0)

(548, 32) (608, 14)

(612, 11)

(620, 46)

(632, 7)

(644, 12)

(644, 27)

(664, 28) (667, 22)

(691, 34)

(697, 9)

Dense array

0. 0.]]

(722, 16)

print('Dense array\n')

output=transform(corpus1, vocab2)

print('Normalised TF\_IDF values')

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0 1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

1.0

print(output[0].toarray(),'\n')

0.7071067811865475 0.7071067811865475

0.7071067811865475

0.7071067811865475

In [78]: #Dense matrix representation for the first document in the corpus

print('Shape of dense matrix',(output[0].shape))

0.37796447300922725
0.37796447300922725

0.37796447300922725

0.37796447300922725

0.37796447300922725
0.37796447300922725

0.37796447300922725

print(output.shape, '\n')

Normalised TF\_IDF values

print("you need to pass list of strings")