```
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) = \frac{\text{Number of times term t appears in a document}}{\text{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:
                $IDF(t) = \log_{e}\frac{\text{Total number of documents}} {\text{Number of documents with term t in it}}.$ for numerical stability
                we will be changing this formula little bit DF(t) = \log_{e}\frac{\text{text}Total number of documents}}{\text{Number of documents}}
                with term t 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.
             </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{text} }{1\text{text} }
                    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
                  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
                  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
                  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 [181]: ## 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
In [168]:
            vectorizer = TfidfVectorizer()
            vectorizer.fit(corpus)
            skl_output = vectorizer.transform(corpus)
In [169]: # 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 [170]: # 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 1.22314355 1.51082562 1.
                                                                  1.91629073 1.91629073
                          1.91629073 1.
In [171]: # shape of sklearn tfidf vectorizer output after applying transform method.
            skl_output.shape
Out[171]: (4, 9)
In [172]: # 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
                               0.38408524091481483
               (0, 6)
               (0, 3)
                               0.38408524091481483
               (0, 2)
                               0.5802858236844359
               (0, 1)
                               0.46979138557992045
In [173]: # 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.
                                                                                Θ.
                                        0.38408524]]
              0.38408524 0.
            Your custom implementation
In [175]: # 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
            from tqdm import tqdm
            from scipy.sparse import csr_matrix
            import math
            import operator
            from sklearn.preprocessing import normalize
            import numpy
            def fit(dataset):
                 unique_words = set() # at first we will initialize an empty set
                 # check if its list type or not
                 if isinstance(dataset, (list,)):
                      for row in dataset: # for each review in the dataset
                          for word in row.split(" "): # for each word in the review. #split method convert
            s a string into list of words
                               if len(word) < 2:
                                    continue
                               unique_words.add(word)
                      unique_words = sorted(list(unique_words))
                      vocab = {j:i for i, j in enumerate(unique_words)}
                      N = len(dataset) #length of dataset
                      # creating empty dict.
                      idf_vals = dict()
                      for i in vocab.keys():
                          idf_vals[i] = 0
                      # c=count
                      C=0
                      # calculate in how many row a world present
                      for i in idf_vals.keys():
                          for doc in dataset:
                               if i in doc:
                                    C += 1
                          idf_vals[i] =c
                          C=0
                      #calculate idf
                      for term, vals in idf_vals.items():
                          idf_vals[term]=1+math.log((1+N)/float(vals+1))
                      return vocab,idf_vals
                 else:
                      print("you need to pass list of strings")
In [179]: vocab,idf_vals=fit(corpus)
In [177]: # nunber of unique word
            # or number of features
            print(vocab)
            {'and': 0, 'document': 1, 'first': 2, 'is': 3, 'one': 4, 'second': 5, 'the': 6, 'third': 7,
            'this': 8}
In [178]: # printing idf value of tf-idf .
            print(list(idf_vals.values()))
            [1.916290731874155, 1.2231435513142097, 1.5108256237659907, 1.0, 1.916290731874155, 1.9162907
            31874155, 1.0, 1.916290731874155, 1.0]
            # now we will use transfrom function to compute tf-idf.
  In [ ]:
In [157]: def transform(dataset, vocab, idf_vals):
                 #Initialiseing empty sparse matrix of size row = length of dataset columns = length of v
            ocab
                 sparse_matrix = csr_matrix((len(dataset),len(vocab)))
                 # in each row in dataset
                 for row in range(0,len(dataset)):
                      #calculating number of words in row and storing in dict.
                      num_of_word_in_row= dict(Counter(dataset[row].split()))
                      for word in dataset[row].split():
                           # if in that row the word present in the vocab.
                          if word in list(vocab.keys()):
                               # implementing tf-idf std_formula .
                               tf_idf=(num_of_word_in_row[word]/len(dataset[row].split()))*(idf_vals[word])
                               # assigning tf_idf to that particular row and columm.
                               sparse_matrix[row, vocab[word]] = tf_idf
                 #here we are applying 12 Normalixation.
                 output =normalize(sparse_matrix, norm='12', axis=1, copy=True, return_norm=False)
                 return output
In [183]: # calling both function
            vocab,idf_vals=fit(corpus)
            output=(transform(corpus, vocab, idf_vals))
In [184]: # Here the output is a sparse matrix of first line of dataset.
            print(output[0])
                               0.4697913855799205
               (0, 1)
               (0, 2)
                               0.580285823684436
               (0, 3)
                               0.3840852409148149
               (0, 6)
                               0.3840852409148149
               (0, 8)
                               0.3840852409148149
In [182]: # first row of tf-idf matrix in dense form
            print(output[0].toarray())
                            0.46979139 0.58028582 0.38408524 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 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 [185]: # Below is the code to load the cleaned_strings pickle file provided
            # Here corpus is of list type
            import pickle
            from operator import itemgetter
            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 [186]: # Write your code here.
            # Try not to hardcode any values.
            # Make sure its well documented and readble with appropriate comments.
In [187]: # apply fit function
            def fit(dataset):
                 unique_words = set() # at first we will initialize an empty set
                 # check if its list type or not
                 if isinstance(dataset, (list,)):
                      for row in dataset: # for each review in the dataset
                          for word in row.split(" "): # for each word in the review. #split method convert
            s a string into list of words
                               if len(word) < 2:</pre>
                                    continue
                               unique_words.add(word)
                      unique_words = sorted(list(unique_words))
                      vocab = {j:i for i,j in enumerate(unique_words)}
                      N = len(dataset) # length of the dataset.
                      # creating empty dict.
                      idf_vals = dict()
                      for i in vocab.keys():
                          idf_vals[i] = 0
                      \# c = count
                      # calculate in how many row a world present
                      for i in idf_vals.keys():
                          for doc in dataset:
                               if i in doc:
                                    c += 1
                          idf_vals[i] =c
                          C=0
                       #calculate idf
                      for term, vals in idf_vals.items():
                          idf_vals[term]=1+math.log((1+N)/float(vals+1))
                      return vocab,idf_vals
                 else:
                      print("you need to pass list of strings")
In [188]: vocab,idf_vals=fit(corpus)
In [189]: # we will have to sort idf values and slice top 50 of them.
In [190]: # sorting in reverse and slice top 50.
            sorted_50_idf_vals = dict(sorted(idf_vals.items(), key=lambda x: x[1],reverse=True)[:50])
            # here we have top 50 idf.
            print(sorted_50_idf_vals)
            {'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, 'accurately': 6.922918004572872, 'achille': 6.922918004572872, 'ackerman': 6.9229180045728
            72, 'adams': 6.922918004572872, 'added': 6.922918004572872, 'admins': 6.922918004572872, 'adm
            iration': 6.922918004572872, 'admitted': 6.922918004572872, 'adrift': 6.922918004572872, 'adv
            enture': 6.922918004572872, 'aesthetically': 6.922918004572872, 'affected': 6.92291800457287
            2, 'affleck': 6.922918004572872, 'afternoon': 6.922918004572872, 'agreed': 6.922918004572872,
            'aimless': 6.922918004572872, 'aired': 6.922918004572872, 'akasha': 6.922918004572872, 'aler
            t': 6.922918004572872, 'alike': 6.922918004572872, 'allison': 6.922918004572872, 'allowing':
            6.922918004572872, 'alongside': 6.922918004572872, 'amateurish': 6.922918004572872, 'amazed':
            6.922918004572872, 'amazingly': 6.922918004572872, 'amusing': 6.922918004572872, 'amust': 6.9
            22918004572872, 'anatomist': 6.922918004572872, 'angela': 6.922918004572872, 'angelina': 6.92
            2918004572872, 'angry': 6.922918004572872, 'anguish': 6.922918004572872, 'angus': 6.922918004
            572872, 'animals': 6.922918004572872, 'animated': 6.922918004572872, 'anita': 6.9229180045728
            72, 'anniversary': 6.922918004572872, 'anthony': 6.922918004572872, 'antithesis': 6.922918004
            572872, 'anyway': 6.922918004572872}
In [191]: # now we have find vocab which have top 50 idf value.
            # all key of top 50 list in new.
            new=list(sorted_50_idf_vals.keys())
            # new_vocab is all features.
            new_vocab = {j:i for i,j in enumerate(new)}
            print(new_vocab)
```

{'aailiyah': 0, 'abandoned': 1, 'abroad': 2, 'abstruse': 3, 'academy': 4, 'accents': 5, 'acce ssible': 6, 'acclaimed': 7, 'accolades': 8, 'accurately': 9, 'achille': 10, 'ackerman': 11, 'adams': 12, 'added': 13, 'admins': 14, 'admiration': 15, 'admitted': 16, 'adrift': 17, 'adve nture': 18, 'aesthetically': 19, 'affected': 20, 'affleck': 21, 'afternoon': 22, 'agreed': 23, 'aimless': 24, 'aired': 25, 'akasha': 26, 'alert': 27, 'alike': 28, 'allison': 29, 'allowing': 30, 'alongside': 31, 'amateurish': 32, 'amazed': 33, 'amazingly': 34, 'amusing': 35, 'amust': 36, 'anatomist': 37, 'angela': 38, 'angelina': 39, 'angry': 40, 'anguish': 41, 'angus': 42, 'animals': 43, 'animated': 44, 'anita': 45, 'anniversary': 46, 'anthony': 47, 'antithesi

#Initialiseing empty sparse matrix of size row = length of dataset columns = length of v

s': 48, 'anyway': 49}

ocab

In [192]: # now we will apply transfer function

def transform(dataset, vocab, idf\_vals):

for row in range(0,len(dataset)):

for word in dataset[row].split():

# in each row in dataset

sparse\_matrix = csr\_matrix((len(dataset),len(vocab)))

#calculating number of words in row and storing in dict.
num\_of\_word\_in\_row= dict(Counter(dataset[row].split()))

# if in that row the word present in the vocab.