r (What does tf-idf mean? If-idf stands for term frequency-inverse document frequency, and the tf-idf weight is a weight often used in information retrieval and te 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 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 function are variants of this simple model. One of the 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
v (Typically, the tf-idf weight is composed by two terms: the first computes the normalized Term Frequency (TF), aka. the number of times 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 the document length (aka. the total number of terms in the document) as a way of normalization: $TF(t) = \frac{\text{Number of times term t appears in a document}}{\text{Total number of terms in the document}}.$ • 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: $IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents with term t in it}}. \text{ for numerical stability we will be changing this formula little bit} \\ IDF(t) = \log_e \frac{\text{Total number of documents}}{\text{Number of documents}}.$
(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
<	Trequency (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 = 0.16 $
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
	 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 ze 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 sklear
	 tfidf vectorizer. Print out the idf values from your implementation and check if its the same as that of sklearns tfidf vectorizer idf values. Once you get your voacb and idf values to be same as that of sklearns implementation of tfidf vectorizer, proceed to the below steps. Make sure the output of your implementation is a sparse matrix. Before generating the final output, you need to normalize you sparse matrix using L2 normalization. You can refer to this link https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html After completing the above steps, print the output of your custom implementation and compare it with sklearns implementation of tfidf vectorizer. To check the output of a single document in your collection of documents, you can convert the sparse matrix related only to the 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 punctuation
v I a	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 your finally submitting the assignment, make sure your code is readable and try not to print things which are not part of this task.
:	<pre>## SkLearn# Collection of string documents Corpus = ['this is the first document', 'this document is the second document',</pre>
	'and this is the third one', 'is this the first document',] SkLearn Implementation
:	<pre>from sklearn.feature_extraction.text import TfidfVectorizer vectorizer = TfidfVectorizer() vectorizer.fit(Corpus) skl_output = vectorizer.transform(Corpus)</pre>
:	<pre># 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']</pre>
:	<pre># 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_)</pre>
:	[1.91629073 1.22314355 1.51082562 1.
:	(4, 9) # sklearn tfidf values for first line of the above corpus. # Here the output is a sparse matrix
	print(skl_output[0]) (0, 8)
:	(0, 1) 0.46979138557992045 # 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 pri # Notice that this output is normalized using L2 normalization. sklearn does this by default.
	print(skl_output.toarray()) [[0.
•	[0.51184851 0.
:	# 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
	<pre>from collections import Counter from tqdm import tqdm from time import sleep from scipy.sparse import csr_matrix import math import operator</pre>
	<pre>from sklearn.preprocessing import normalize import numpy as np import os import warnings warnings.filterwarnings("ignore") import pandas as pd import math</pre>
:	<pre>def fit(corpus): word_dict = set() # set the word_dict as empty set if isinstance(corpus,(list,)):# check for the corpus type for row in corpus:</pre>
	<pre>for word in row.split(" "): #slpit the each word in the row if len(word) < 2: # remove the word with length less than 2 continue word_dict.add(word) # add the word into the dict word_dict = sorted(list(word_dict)) # sort the dict and convert in into the list vocab = {i:j for j,i in enumerate(word_dict)}# create the dict of word with their index value</pre>
	<pre>doc_word_freq=[] idf_value=[] for word in word_dict: doc_num = 0 for doc in corpus: if word in doc.split():</pre>
	<pre>doc_num doc_word_freq.append(doc_num) for i in range(len(doc_word_freq)): idf = 1+math.log((1+len(corpus))/(1+(doc_word_freq[i]))) idf_value.append(idf) return(vocab,idf_value) else:</pre>
	<pre>print("Please provide the data in list") def transform(corpus, vocab, idf_value): rows = [] columns = [] values = [] tf_value=[]</pre>
	<pre>tf_value=[] tf_idf_value = [] Sparse_Matrix=csr_matrix((len(corpus),len(vocab))) if isinstance(corpus, (list,)): for idx, row in enumerate(tqdm(corpus)):# for each row in the corpus</pre>
	<pre>for word, freq in word_freq.items():# for each unique word in the review. if len(word) < 2: continue # we will check if its there in the vocabulary that we build in fit() function # dict.get() function return the values, if key doesn't exits it return -1 col index = vocab.get(word, -1)</pre>
	<pre>if col_index !=-1: rows.append(idx) #storing, index of the document columns.append(col_index) #storing, dimensions of the word values.append(freq/len(row)) #storing, frequency of the w r # for doc in corpus: # tf=[]</pre>
	# doc = doc.split(' ') # for word in vocab.keys(): # t_freq = 0 # for words in range(len(doc)): # if word == doc[words]: # t_freq += 1
	<pre>#</pre>
:	<pre># return normalize(Sparse_Matrix,norm='12') else: print("Please provide the data in list") def fit(corpus):</pre>
	<pre>word_dict = set() # set the word_dict as empty set if isinstance(corpus,(list,)):# check for the corpus type for row in corpus: for word in row.split(" "): #slpit the each word in the row if len(word) < 2: # remove the word with length less than 2</pre>
	<pre>word_dict = sorted(list(word_dict)) # sort the dict and convert in into the list vocab = {i:j for j,i in enumerate(word_dict)}# create the dict of word with their index value doc_word_freq=[] idf_value=[] for word in word_dict: doc_num = 0</pre>
	<pre>for doc in corpus: if word in doc.split(): doc_num +=1 doc_word_freq.append(doc_num) for i in range(len(doc_word_freq)): idf = 1+math.log((1+len(corpus))/(1+(doc_word_freq[i]))) idf value.append(idf)</pre>
	<pre>return(vocab, idf_value) else: print("Please provide the data in list") def transform(corpus, vocab, idf_value): rows = []</pre>
	<pre>columns = [] values = [] tf_value=[] tf_idf_value = [] Sparse_Matrix=csr_matrix((len(corpus),len(vocab))) if isinstance(corpus, (list,)):</pre>
	<pre>for idx, row in enumerate(tqdm(corpus)):# for each row in the corpus # it will return a dict where key is the word and freq as a values, for every unique word in t word_freq = Counter(row.split()) for word, freq in word_freq.items():# for each unique word in the review. if len(word) < 2: continue # we will check if its there in the vocabulary that we build in fit() function</pre>
	<pre># dict.get() function return the values, if key doesn't exits it return -1 col_index = vocab.get(word, -1) if col_index !=-1: idf = idf_value[col_index]</pre>
	<pre>rows.append(idx) #storing, index of the document columns.append(col_index) #storing, dimensions of the word</pre>
:	<pre>columns.append(col_index) #storing, dimensions of the word</pre>
:	<pre>columns.append(col_index) #storing, dimensions of the word</pre>
	<pre>columns.append(col_index) #storing, dimensions of the word</pre>
	<pre>columns.append(col_index) #storing, dimensions of the word</pre>
:	<pre>columns.append(col_index) #storing, dimensions of the word tf_idf_value.append((freq/len(row))*idf) #storing, frequency of the w return normalize(csr_matrix((tf_idf_value, (rows,columns)), shape=(len(corpus),len(vocab))),norm='12') Corpus = ["this is the first document", "this document is the second document", "and this is the third one", "is this the first document",] vacob,idf_value = fit(Corpus) print("using the fit() idf values:\n",idf_value) print("using the fit() vocab:\n",vacob) using the fit() idf values: [1.916290731874155, 1.2231435513142097, 1.5108256237659907, 1.0, 1.916290731874155, 1.916290731874155, 1.0] using the fit() vocab: ('and': 0, 'document': 1, 'first': 2, 'is': 3, 'one': 4, 'second': 5, 'the': 6, 'third': 7, 'this': 8} from sys import getsizeof # print(getsizeof(transform(Corpus, vacob,idf_value))) print(transform(Corpus,vacob,idf_value)[0])</pre>
:	<pre>columns.append(col_index) #storing, dimensions of the word</pre>
	columns.append(col_index) **storing, dimensions of the word
	columns.append(col index) **storing, dimensions of the word tfidf_value.append((freq/len(row))*idf) **storing, frequency of the w return normalize(csr_matrix({tf_idf_value, (rows,columns)}, shape=(len(corpus),len(vocab))),norm='12') Corpus = ["this is the first document", "this document is the second document", "and this is the third one", "le this the first document", "le this the first document", "vising the fit() idf values:\n",idf_value) print("using the fit() idf values:\n",vacob) using the fit() idf values: [1,916290731874155, 1.2231435513142097, 1.5108256237659907, 1.0, 1.916290731874155, 1.916290731874155, 1.0 sing the fit() vocab: ('and': 0, 'document': 1, 'first': 2, 'is': 3, 'one': 4, 'second': 5, 'the': 6, 'third': 7, 'this': 8} from sys import getsizeof *print(getsizeof(transform(Corpus, vacob,idf_value))) print(transform(Corpus,vacob,idf_value)[0]) 1003[(0, 1)
	columns.append(col_index) **Storing, dimensions of the word tf_idf_value.append(tfseq_len(con)*idf) **Storing, frequency of the w return normalize(csr_matrix((tf_idf_value, (rows,columns)), shape=(len(corpus),len(vocab))),norm='12') Corpus = { "this is the first document", "this document is the second document", "and this is the third one", "is this the first document", "y vacob_idf_value = fit(Corpus) print("using the fit() idf values:\n",idf_value) print("using the fit() vocab:\n",vacob) using the fit() idf values: [1.916290731874155, 1.2231435513142097, 1.5108256237659907, 1.0, 1.916290731874155, 1.916290731874
	columns.append(col_index) **actoring, dimensions of the word return normalize(car_matrix((tr_idf_value, (rows,columns)), shape=(len(corpus),len(vocab))),norm='12'; Corpus = { "this is the first document", "mand this is the first document", "and this is the first document", "so this the first document', "so this first the first doc
	columns.appendioco_indexy' secoring, disensions of the word tt_iof_value.appendictors/locations/light/sactions/requirecy of the w return normalize(ins_notrix((inf_inf_value, irous,columns)), shape=(len(sorpusi,len(woos))),norm='12'; Cnepus = ' "init document is the second document", "init document is the second document", "and this is the first document", "so this fert document", "uncomplete the fit() (orgus) vacob,idf_value = fit(Corpus) print("second the fit() woosb\n",vacob) using the fit() idf values: [1.916920731874155, 1.232434553142097, 1.510825623765907, 1.0, 1.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155, 1.316330731874155, 2.916290731874155
	return normalizates; matrix(tif_iof_value, (rows_columns)), shape=[lentcopyas], lentvocab)), norm="12"; Coppus = { "this is the first document", "and this the first document", "this is the first document", "this is the list of content on", "and this the first document", "this is the fir
	columnic append (circl_indexs) * Secretary, Statematics of the word to its ist values.ppend (circl_value, lrows.oolumnal), shape=(lentoorpus).lentwool)!,norm='12'; **Cerpus = { **Cerpus = { **Cerpus = { **Tatis is the first document*; **Tatis document* is the second document*, **Tatis document list the second document*, **Tatis document list the second document*, **Tatis document* is the second document*,
	columns appeared real jurdent enteriors, differences one of the word relative appeared real jurdent, reported realization, foregouncy of the word relative products and the control of the company of the word relative to the company of the compa
	return normalize (or_maintenance) the vector according to the vector of
	return normalization_processing_color_colo
	concursi approaches. Instead * **Charley**, immunities for swell to first and supposed(toog) factors (**Tables) (**Color additional)**, anapose (**Tables) (**Ta
	Solume's approaches, index of sections of control of the control o
	included appearation (colors) settlement and the send of the send
	returns accompliant in approximately placed in a function of communication of the control of communication o
	Copye = 1 The property of the
	Columbia grantered placed in columbia of the c
	Course = {
	Deligned a process of the control process of

print(transform(corpus, vocab,top_idf)[0].toarray())

| 746/746 [00:00<00:00, 7479