collection or corpuengines as a central	idf mean? erm 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 us. 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 all tool in scoring and ranking a document's relevance given a user query. est ranking functions is computed by summing the tf-idf for each query term; many more sophisticated ranking functions are variants of this simple model. essfully used for stop-words filtering in various subject fields including text summarization and classification.
Inverse Document • TF: Term From frequency is $TF(t) = \frac{\text{Num}}{t}$ • IDF: Inverse have little im $IDF(t) = \log t$ Example Consider a document	If 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 to the 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. **Requency, 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 often divided by the document length (aka. the total number of terms in the document) as a way of normalization: **More of times term t appears in a document.** **Total number of terms in the document.** **Pocument 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 apportance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following: **Total number of documents** **Total number of docu
 As a part of t You should c Sklearn does Sklearn Sklearn which p 	The Vectorizer & compare its results with Sklearn: this task you will be implementing TFIDF vectorizer on a collection of text documents. The properties of your own implementation of TFIDF vectorizer with that of sklearns implementation TFIDF vectorizer. The word 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: The initial properties are the vectorizer of the initial properties are the vectorizer of the idea of the initial properties are of visions. $IDF(t) = 1 + \log_e \frac{1 + Total number of documents in collection}{1 + Number of documents with term t in it.}$
4. The fina • Steps to appr 1. You won 2. Print ou 3. Print ou 4. Once you 5. Make sure learn.org 6. After co 7. To check	applies L2-normalization on its output matrix. all output of sklearn tfidf vectorizer is a sparse matrix. coach this task: uld have to write both fit and transform methods for your custom implementation of tfidf vectorizer. It 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. It the idf values from your implementation and check if its the same as that of sklearns fidf vectorizer idf values. For get your voacb and idf values to be same as that of sklearns implementation of tfidf vectorizer, proceed to the below steps. For each 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-g/stable/modules/generated/sklearn.preprocessing.normalize.html For pulleting the above steps, print the output of your custom implementation and compare it with sklearns implementation of tfidf vectorizer. It is the fidf 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 letters or punctuat Note-3: During the not part of this tas Corpus ## SkLearn# Corpus = ['theory and thie 'is this import warning	aut 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 tions, 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 is 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 ask. **Collection of string documents** nis is the first document', by the first document', stip is the first document', stip is the first document']
<pre>from collecti from tqdm imp from scipy.sp import math import operat from sklearn. import numpy corpus = ['th 'thi 'and 'is def fit(doc): """given This f if isinst</pre>	cons import Counter cort tqdm coarse import csr_matrix cor cor preprocessing import normalize as np mis is the first document', its document is t the second document', it this is the third one', this the first document']
words words retur else: print #comparison vocab = fit(comparist) from sklearn. # sklearn feat print("sklear	<pre>cow in doc: ior i in row.split(): if len(i)<2: continue words.add(i) s=sorted(words) #sorting words in aphabetic order : it helps in enumeration s={v:k for k,v in enumerate(words)} rn words c("Please enter the list of strings") corpus) n:",vocab, '\n') ifeature_extraction.text import TfidfVectorizer ature names,sorts by alphabetic_order as default. n:",TfidfVectorizer().fit(corpus).get_feature_names())</pre>
sklearn: ['an def idf(corpu '''This f idf={} n_docs=le for ind,r for c f	function fids idf for corpus''' en(corpus) #number of documents in corpus row in enumerate(corpus): word in vocab.keys(): count=0 for row in corpus: if word in set(row.split()): count+=1 IDF=1+math.log((1+n_docs)/(1+count)) idf[word]=IDF ict(sorted(idf.items(), key=lambda x:x[1], reverse=True)[:50])) iff
290731874155, #comparison print("custom print("sklear custom: [1.91 sklearn: [1.9 1. 1 def transform '''This f print('→ if isinst numbe	2990731874155, 'document': 1.2231435513142097, 'first': 1.5108256237659907, 'is': 1.0, 'one': 1.916290731874155, 'second': 1.916290731874155, 'the': 1.0, 'third': 'this': 1.0} n:", list(idf(corpus,vocab).values()),"\n") n:", TfidfVectorizer().fit(corpus).idf_) 1.6290731874155, 1.2231435513142097, 1.5108256237659907, 1.0, 1.916290731874155, 1.916290731874155, 1.0, 1.916290731874155, 1.0] 1.01629073 1.22314355 1.51082562 1.
#In h frequence print #2 csr=c #help	<pre># [row.split() for row in corpus] # list of lists of words in each row lency_of_words=[]</pre>
returelse: print #comparison x=transform(cont) print(f'custon from sklearn. vectorizer = vectorizer.fi skl_output = print('\n\nsk	om:\n{x}\nshape:{x.shape}') feature_extraction.text import TfidfVectorizer TfidfVectorizer()
e.g,x=transf print	1.
(2, 8) (3, 1) (3, 2) (3, 3) (3, 6) (3, 8) shape:(4, 9) sklearn: (0, 8) (0, 6) (0, 3) (0, 2) (0, 1) (1, 8) (1, 6) (1, 5) (1, 3) (1, 1)	0.267103787642168 0.4697913855799205 0.58025823684436 0.3840852409148149 0.3840852409148149 0.3840852409148149 0.38408524091481483 0.38408524091481483 0.38408524091481483 0.38408524091481483 0.5802658236844359 0.46979138557992045 0.281088674033753 0.281088674033753 0.281088674033753 0.281088674033753 0.6876235979836938
_	max features functionality:
 This task is s of documents Here you wil Steps to appr You won Now son your von Make sun learn.org 	uld 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 rt 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
<pre>#Here corpus import pickle with open('cl corpus = #printing the print("Number Number of doc def fit(doc): """given This f if isinst words</pre>	Leaned_strings', 'rb') as f: pickle.load(f) e length of the corpus loaded r of documents in corpus = ",len(corpus)) cuments in corpus = 746
words words retur else: print vocab = fit(c print(vocab)) {'aailiyah': 11, 'accurate 3, 'actors': 35, 'adorable	for i in row.split(): if len(i)>=2: words.add(i) else: continue s=sorted(words) #sorting words in aphabetic order : it helps in enumeration s={v:k for k,v in enumerate(words)} *n words c("Please enter the list of strings") corpus) 0, 'abandoned': 1, 'ability': 2, 'abroad': 3, 'absolutely': 4, 'abstruse': 5, 'abysmal': 6, 'academy': 7, 'accents': 8, 'accessible': 9, 'acclaimed': 10, 'accolaite': 12, 'accurately': 13, 'accused': 14, 'achievement': 15, 'achille': 16, 'ackerman': 17, 'act': 18, 'acted': 19, 'acting': 20, 'action': 21, 'actions': 22, 'act 24, 'actress': 25, 'actresses': 26, 'actually': 27, 'adams': 28, 'adaptation': 29, 'add': 30, 'added': 31, 'addition': 32, 'admins': 33, 'admiration': 34, 'admint': 37, 'adventure': 38, 'advise': 39, 'aerial': 40, 'aesthetically': 41, 'affected': 42, 'affleck': 43, 'afraid': 44, 'africa': 45, 'afternoon': 46 '!': 48, 'ages': 49, 'ago': 50, 'agree': 51, 'agreed': 52, 'aimless': 53, 'air': 54, 'aired': 55, 'akasha': 56, 'akin': 57, 'alert': 58, 'alexander': 59, 'alike': 58, 'alexander': 59, 'alike': 58, 'alexander': 59, 'alike': 58, 'alexander': 58, 'alexa
rocity': 142, 'austere': 15 4, 'awkwardly 176, 'bag': 1 88, 'barney': c': 200, 'bec 'believe': 21 t': 223, 'bet chy': 235, 'b 'boasts': 248 'boogeyman': 1, 'bother': revity': 283, e': 294, 'bro 'bullock': 30 'came': 318, ptain': 330, 1, 'cars': 34 ses': 354, 'c ard': 365, 'c	132, 'aspects': 133, 'ass': 134, 'assante': 135, 'assaulted': 136, 'assistant': 137, 'astonishingly': 138, 'astronaut': 139, 'atmosphere': 140, 'atrocious': 14 'attempti': 143, 'attempted': 144, 'attemptind': 145, 'attempts': 146, 'attention': 147, 'attractive': 148, 'audience': 149, 'audio': 150, 'aurv': 151, 'austen': 133, 'author': 154, 'average': 155, 'aversion': 156, 'avoid': 157, 'avoided': 158, 'award': 159, 'awarded': 160, 'awards': 161, 'away': 162, 'awesome': 163, 'awful '156, 'aye': 166, 'baaaaaad': 167, 'babbling': 168, 'bable': 169, 'baby': 170, 'bablet': 183, 'balls': 184, 'band': 173, 'backd': 174, 'bade': 175, 'backl': 172, 'backdrop': 173, 'backd': 174, 'bade': 175, 'bashe': 189, 'barren': 190, 'based': 191, 'basic': 192, 'basically': 193, 'bat': 194, 'bates': 195, 'baxendale': 196, 'bear': 197, 'beautiful': 198, 'beautifully': 199, 'ame': 201, 'bechard': 202, 'become': 203, 'becomes': 204, 'began': 205, 'begin': 206, 'beginning': 207, 'behind': 208, 'behold': 209, 'bela': 210, 'believable': 21, 'believed': 213, 'bell': 214, 'bellucci': 215, 'belly': 216, 'belmod': 217, 'ben': 218, 'bending': 219, 'bennett': 220, 'bergen': 221, 'bertolucci': 222, 'ter': 224, 'betty': 225, 'beware': 226, 'beyond': 227, 'bible': 228, 'big': 229, 'biggest': 230, 'billy': 231, 'biographical': 232, 'bipolarity': 233, 'bit': 234, 'black': 236, 'bland': 237, 'blake': 238, 'bland': 239, 'blandly': 240, 'blare': 241, 'blatant': 242, 'blew': 243, 'blood': 244, 'blown': 257, 'bonuses': 258, 'bonbs': 260, 'book': 261, 'boost': 262, 'bop': 263, 'bordered': 264, 'borderlines': 266, 'borders': 266, 'border': 266, 'border': 276, 'bose': 267, 'bored': 268, 'boring': 269, 'boring': 269, 'borrowed': 270, 'bose': 272, 'bothersome': 273, 'bought': 274, 'box': 275, 'boyfriend': 276, 'boyle': 277, 'brain': 278, 'brainsucking': 279, 'brait': 280, 'breaking': 281, 'brian': 282, 'brian': 286, 'borther': 286, 'brorther': 286, 'burder': 381, 'brailly': 307, 'bunch': 388, 'burton': 389, 'burton': 389, 'burton': 389, 'burton': 389, 'burto
n': 397, 'chi 'christopher' 418, 'clear': 0, 'co': 431, 42, 'colours' 'comment': 45 mpleted': 464 473, 'concept nfirm': 483, le': 493, 'co 2, 'containin t': 512, 'con oluted': 522, 534, 'couple' d': 546, 'cra 557, 'credits 9, 'cruise': ds': 582, 'da 594, 'dead': 'decidely': 6 'delight': 61	(ldrens': 388, 'cheerfull': 389, 'cheerless': 390, 'cheesness': 391, 'cheesy': 392, 'chemistry': 393, 'chick': 394, 'child': 395, 'childhood': 396, 'childrens': 398, 'chills': 399, 'chilly': 400, 'chimp': 401, 'chodorov': 402, 'choice': 403, 'choices': 404, 'choked': 405, 'chosen': 406, 'chow': 407, 'christmas': 409, 'church': 410, 'cinemat': 411, 'cinematic': 412, 'cinematographers': 413, 'cinematography': 414, 'circumstances': 415, 'class': 416, 'classic': 417, 'class 419, 'clearly': 420, 'clever': 421, 'clich': 422, 'cliche': 423, 'collerts': 424, 'cliff': 425, 'climax': 426, 'colse': 427, 'closed': 428, 'clothes': 429, 'clube': 431, 'collert': 432, 'collect': 439, 'collective': 440, 'coored': 441, 'colorf': 443, 'collembe': 444, 'come': 445, 'comedic': 446, 'comedy': 447, 'comes': 448, 'comfortable': 449, 'comforting': 459, 'comical': 451, 'coming': 452, 'commands': 447, 'commentary': 455, 'commented': 456, 'complexity': 467, 'complexity': 467, 'composed': 468, 'confoses': 497, 'consestentle': 470, 'composed': 468, 'confoses': 489, 'confidence': 480, 'confidence': 481, 'confidence': 481, 'confidence': 481, 'confidence': 481, 'confidence': 482, 'confidence': 484, 'confuses': 485, 'confusing': 486, 'connections': 487, 'connections': 489, 'considering': 490, 'considering': 491, 'considering': 494, 'considering': 495, 'considers': 496, 'considens': 487, 'consolations': 489, 'constatine': 590, 'constructed': 591, 'contained ag': 593, 'contains': 594, 'contributory': 514, 'contrived': 515, 'control': 516, 'control': 517, 'continue': 529, 'continuity': 599, 'continuously': 519, 'contract': 511, 'contributing': 513, 'contributory': 514, 'control': 526, 'corn': 527, 'conversy': 527, 'conversy': 528, 'core': 529, 'const': 529, 'const': 529, 'const': 520, 'const': 52
p': 649, 'dev ferent': 660, ointed': 670, g': 679, 'dis 89, 'diving': 'doomed': 701 713, 'drawn': 'duet': 726, h': 738, 'eas d': 750, 'eff 61, 'elegant' o': 772, 'eming': 783, 'en 4, 'enjoy': 7 ly': 805, 'en nce': 816, 'e 'everybody': e': 838, 'excuciatingly':	10, 'designed': 640, 'designer': 641, 'desperately': 642, 'desperation': 643, 'despised': 644, 'despite': 645, 'desperoy:: 646, 'detailing': 647, 'details': 648, 'deelopment': 650, 'developments': 651, 'di': 652, 'diabetic': 653, 'dialog': 654, 'dialogue': 656, 'diaper': 657, 'dickens': 658, 'difference': 659, 'digner': 667, 'disperi: 667, 'disperi: 667, 'disperi: 667, 'disperi: 667, 'disperi: 667, 'disperi: 668, 'director': 668, 'director': 669, 'director': 667, 'disperi: 667, 'disperiting': 671, 'disperiting': 672, 'disperiting': 673, 'disperiting': 674, 'disperiting': 674, 'disperiting': 674, 'disperiting': 675, 'disperiting': 676, 'disperiting': 677, 'disprace': 678, 'disperiting': 680, 'doctor': 681, 'documentaries': 692, 'documentary: 693, 'dodge': 694, 'dogs': 695, 'dollars': 696, 'dominated': 687, 'done': 688, 'disturbi 690, 'doctor': 691, 'documentaries': 692, 'documentary: 693, 'dodge': 694, 'dogs': 695, 'dollars': 696, 'dominated': 687, 'done': 698, 'donlevy: 699, 'dont': -, 'dose': 702, 'doubt': 703, 'downs': 704, 'dozen': 705, 'dr': 706, 'dracula': 707, 'draft': 708, 'drag': 709, 'drago': 710, 'drama': 711, 'dramatic': 712, 'draw '714, 'dream': 715, 'dreams': 716, 'dreary': 717, 'dribble': 718, 'drifting': 720, 'drive': 721, 'drobing': 722, 'dropped': 723, 'dry': 724, 'due' 'dull': 727, 'dumb': 728, 'dumbest': 729, 'duper': 730, 'duris': 731, 'dustin': 732, 'dvd': 733, 'dwight': 734, 'dysfunction': 735, 'earlier': 736, 'early': 737, 'diy': 739, 'easy': 740, 'eating': 741, 'ebay': 742, 'ebola': 743, 'eccleston': 744, 'ed': 745, 'edge': 746, 'editing': 747, 'edition': 748, 'edct': 751, 'effective': 752, 'effectis': 753, 'effort': 754, 'efforts': 755, 'egotism': 756, 'eighth': 757, 'eimbir': 758, 'einber': 759, 'enbarsaly': 766, 'elsewhere': 767, 'embarrassing': 768, 'embarrassing': 769, 'embassy': 770, 'emerge': 771, 'et': 771, 'embarrassing': 796, 'enjoyed': 797, 'english': 798, 'enotion': 778, 'enotion': 778, 'enotion': 778, 'enotion': 778, 'enotion': 778, 'enotion': 778, 'enotion': 798, 'enotion':
8, 'face': 87 ling': 891, ' asy': 903, 'f 4, 'faultless t': 926, 'fei 8, 'film': 93 ds': 950, 'fi 62, 'flakes': 974, 'flying' ed': 986, 'fo m': 997, 'for ankly': 1009, 019, 'frontie ntal': 1030, 'games': 1042 'generally': 'gets': 1065, 75, 'give': 1 1087, 'going' 1099, 'govern	expression': 869, 'exquisite': 870, 'extant': 871, 'exteriors': 872, 'extraneous': 873, 'extraordinary': 874, 'extremely': 875, 'eye': 876, 'eyes': 876, 'eyes': 876, 'eyes': 876, 'eyes': 876, 'faces': 880, 'facial': 881, 'facing': 882, 'fact': 883, 'factor': 883, 'factor': 887, 'facing': 887, 'factor': 887, 'facing': 887, 'facing': 888, 'factor': 887, 'facing': 887, 'factor': 887, 'facing': 888, 'factor': 887, 'facing': 888, 'factor': 992, 'ar': 994, 'farce': 995, 'fare': 996, 'fascinated': 997, 'fascinating': 998, 'fascination': 999, 'fashioned': 910, 'fast': 911, 'faster': 912, 'fat': 913, 'father': 917, 'favorite': 918, 'favorite': 919, 'fear': 920, 'feature': 921, 'features': 922, 'feeling': 924, 'feeling': 924, 'feeling': 924, 'feeling': 927, 'fellowes': 928, 'felt': 929, 'female': 930, 'females': 931, 'ferry': 932, 'fest': 933, 'fi': 934, 'fields': 935, 'fifteen': 936, 'fifties': 937, 'filmed': 940, 'filming': 941, 'filmed': 942, 'filmography': 943, 'films': 944, 'final': 945, 'finale': 946, 'finally': 947, 'financial': 948, 'find': 954, 'fire': 952, 'fingernails': 953, 'finished': 954, 'fire': 955, 'fish: 957, 'fishnet': 952, 'fingernails': 953, 'finished': 954, 'fire': 955, 'fish: 957, 'fishnet': 958, 'fisted': 959, 'fit': 960, 'five': 961, 'fle': 963, 'fle': 964, 'fle': 957, 'footage': 977, 'fodder': 978, 'follow': 979, 'following': 988, 'forest': 979, 'footage': 983, 'football': 984, 'force': 985, 'force': 987, 'ford': 988, 'foregn': 989, 'foreigner': 999, 'forever': 991, 'forget': 992, 'forgettable': 993, 'forgetting': 994, 'forgot': 995, 'forgetten': 996, 'face': 995, 'former': 999, 'fort': 1000, 'forth': 1001, 'forwarded': 1002, 'found': 1003, 'four': 1004, 'foot': 1005, 'footage': 983, 'football': 984, 'force': 995, 'found': 1004, 'foot': 1004, 'foot': 1005, 'footage': 994, 'forgot': 995, 'found': 1004, 'foot': 1004, 'foot': 1006, 'friendship': 1017, 'frightening': 1018, 'foot': 1027, 'fun': 1028, 'function': 1029, 'foot': 1029, 'footage': 1031, 'ferend': 1033, 'fac': 1033, 'fac': 1034, 'gade': 10
2, 'halfway': ned': 1144, ' 1155, 'hbo': st': 1167, 'h 1178, 'heroes 189, 'hill': s': 1201, 'ho peless': 1212 222, 'hours': s': 1234, 'hu ied': 1245, ' 1255, 'imagin n': 1265, 'im g': 1274, 'in 2, 'indictmen ly': 1292, 'i ead': 1302, ' 1311, 'intens 1320, 'interv	less': 1122, 'guests': 1123, 'guilt': 1124, 'gung': 1125, 'guy': 1126, 'guys': 1127, 'hackneyed': 1128, 'hanggis': 1129, 'hair': 1130, 'hairsplitting': 1131, 'halfe': 1134, 'hand': 1135, 'handle': 1136, 'handles': 1138, 'hands': 1139, 'hang': 1140, 'hankies': 1141, 'hanks': 1142, 'happen': 1143, 'happi': 1146, 'hard': 1147, 'harris': 1148, 'hate': 1149, 'hated': 1150, 'hatred': 1151, 'havilland': 1152, 'hay': 1153, 'hayao': 1154, 'hayben': 1166, 'head': 1157, 'heads': 1158, 'hear': 1159, 'heard': 1160, 'heart': 1161, 'hearts': 1162, 'heartwarming': 1163, 'heaven': 1164, 'heche': 1165, 'heels': 1166, 'hell': 1169, 'hellish': 1170, 'helpis': 1171, 'helpis': 1173, 'helps': 1174, 'hence': 1175, 'hendrikson': 1176, 'hernandez': 1177, 'he': 1179, 'heroine': 1180, 'highest': 1186, 'highlights': 1187, 'highly!: 1188, 'hilario 1190, 'hilt': 1191, 'hip': 1192, 'history': 1193, 'hitchcock': 1194, 'ho': 1195, 'hockey': 1196, 'hoffman': 1197, 'hold': 1198, 'holding': 1199, 'holds': 1200, 'ollander': 1202, 'hollow': 1203, 'hollywood': 1204, 'horrible': 1215, 'horrid': 1216, 'horrified': 1217, 'horror': 1218, 'horse': 1219, 'hosting': 1220, 'hopet': 1210, 'horrendously': 1214, 'horrible': 1215, 'horrid': 1216, 'horrified': 1217, 'horror': 1218, 'horse': 1219, 'hosting': 1220, 'house': 1223, 'house': 1224, 'houses': 1225, 'howdy': 1226, 'howe': 1227, 'howell': 1228, 'however': 1229, 'huge': 1230, 'hugo': 1231, 'humani: 1232, 'humanity': 1233, 'maline': 1256, 'idmb': 1257, 'imitation': 1258, 'imagination': 1254, 'impirial': 1266, 'impirial': 1267, 'improvement': 1268, 'imporvement': 1268, 'improvement': 1269, 'impirial': 1260, 'implausible': 1261, 'impropropriate': 1271, 'incendiary': 1272, 'includes': 1273, 'inconsistencies': 1276, 'incorrectness': 1277, 'incredible': 1278, 'incorpriate': 1279, 'inceperice': 1290, 'inspiring': 1300, 'instruments': 1303, 'insulin': 1304, 'insulin': 1304, 'interpett': 1316, 'interpating': 1308, 'inspiring': 1309, 'install': 1301, 'interving': 1312, 'interions': 1313, 'interpiac': 1314, 'interpet
a': 1353, 'je 'jones': 1366 'junkyard': 1 89, 'kidnappe 'knew': 1401, e': 1412, 'la 'laselva': 14 h': 1435, 'la 6, 'leaves': 1458, 'lets': t': 1470, 'li 81, 'lines': l': 1493, 'lo 'looking': 15 'love': 1516, rics': 1528, 539, 'makes': bles': 1551, 'master': 156 e': 1572, 'mc	12, 'james': 1343, 'jamie': 1344, 'japanese': 1345, 'jason': 1346, 'jay': 1347, 'jealousy': 1348, 'jean': 1349, 'jennifer': 1356, 'jerky': 1351, 'jerry': 1352, 'jsessice': 1354, 'jet': 1355, 'jim': 1356, 'jimmy': 1357, 'job': 1358, 'jobs': 1359, 'joe': 1360, 'john': 1361, 'joins': 1362, 'joke': 1363, 'jokes': 1364, 'jonah': 3, 'jokes': 1364, 'jonah': 1375, 'june': 1376, 'junk': 1376, 'junk': 1376, 'junk': 1375, 'june': 1376, 'junk': 1376, 'junk': 1375, 'june': 1376, 'junk': 1378, 'justice': 1379, 'jutland': 1380, 'kanaly': 1381, 'kathy': 1382, 'keep': 1383, 'keeps': 1384, 'keira': 1385, 'keith': 1386, 'kept': 1387, 'kevin': 1388, 'kind': 1390, 'kids': 1391, 'kieslowski': 1392, 'kill': 1393, 'killer': 1394, 'killing': 1395, 'killings': 1396, 'kind': 1397, 'kinda': 1398, 'kirk': 1399, 'kitchy': 'kinjghtley': 1402, 'knocked': 1403, 'know': 1404, 'known': 1405, 'knows': 1406, 'koteas': 1407, 'kris': 1408, 'kristoffersen': 1409, 'kudos': 1410, 'la': 1411, 'lake': 1413, 'lacked': 1414, 'lackes': 1415, 'ladies': 1416, 'lady': 1417, 'lame': 1418, 'lance': 1419, 'landscapes': 1420, 'lane': 1421, 'lange': 1422, 'largely': 1242, 'laste': 1425, 'last': 1426, 'lasting': 1427, 'latched': 1428, 'late': 1429, 'later': 1430, 'latest': 1431, 'latifa': 1432, 'latin': 1433, 'latin': 1434, 'lagab': 1445, 'lagab
on': 1593, 'm 'microsoft': irrormask': 1 t': 1625, 'mo 5, 'monumenta 1646, 'move': der': 1658, ' 68, 'nasty': d': 1679, 'ne 'never': 1690 'nobody': 170 table': 1713, 724, 'nuts': d': 1735, 'of 1747, 'one': 1758, 'origin eracting': 17 'owns': 1780, 'palance': 17 1803, 'partic	'melodrama': 1584, 'merville': 1585, 'member': 1586, 'members': 1587, 'memorable': 1588, 'memories': 1589, 'memorized': 1599, 'menacic': 1591, 'menacin': 1592, 'mercy': 1594, 'meredith': 1595, 'merit': 1596, 'mesmerising': 1597, 'mess': 1598, 'messages': 1599, 'meteorite': 1600, 'mexican': 1601, 'michael': 1602, 'mickey': 1604, 'middle': 1605, 'might': 1606, 'mighty': 1607, 'mind': 1608, 'mindblowing': 1609, 'miner': 1610, 'mini': 1611, 'minor': 1612, 'miyazaki': 1623, 'minutes': 16.615, 'miserable': 1616, 'miserably': 1617, 'mishima': 1618, 'misplace': 1619, 'miss': 1620, 'missed': 1621, 'mistakes': 1622, 'miyazaki': 1623, 'modern': 1624, 'mollusk': 1626, 'moment': 1627, 'moments': 1638, 'moros': 1639, 'money': 1630, 'monica': 1631, 'monico': 1632, 'montoonous': 1633, 'monster': 1634, 'monstrous' 1612, 'motivations': 1643, 'mountain': 1644, 'mouse': 1645, 'morola': 1643, 'movements': 1649, 'moves': 1659, 'movie': 1651, 'movie': 1652, 'moving': 1653, 'ms': 1654, 'much': 1655, 'muddled': 1656, 'muppets': 1657, 'marrative': 1659, 'murdering': 1660, 'murky': 1661, 'movie': 1662, 'music': 1663, 'must': 1664, 'must': 1664, 'must': 1675, 'nearly': 1676, 'necklace': 1677, 'neartive': 1689, 'newilles': 1681, 'negulesco': 1682, 'neighbour': 1683, 'neil': 1684, 'nerves': 1685, 'nervous': 1686, 'net': 1687, 'neflix': 1688, 'network': 10, 'noreworthy': 1716, 'nothing': 1717, 'novella': 1708, 'normally': 1709, 'northern': 1710, 'nostalgia': 1711, 'not': 1711, 'not': 1712, 'note': 1715, 'noteworthy': 1716, 'nothing': 1717, 'novella': 1718, 'number': 1719, 'numbers': 1720, 'nun': 1721, 'nuns': 1722, 'odd': 1734, 'oder': 1732, 'oder': 1732, 'odd': 1734, 'occur': 1733, 'offers': 1738, 'offers': 1744, 'ooe': 1755, 'oriented': 1756, 'outsa': 1757, 'original 's': 1758, 'original': 1759, 'opening': 1750, 'opening': 1752, 'opening': 1753, 'opinion': 1754, 'ordeal': 1755, 'oriented': 1766, 'outsa': 1767, 'outward': 1766, 'overall': 1770, '
1824, 'perfor g': 1834, 'ph 44, 'picture' lain': 1856, ed': 1867, 'p ignant': 1879 r': 1889, 'po bly': 1899, '909, 'predict sident': 1919 dings': 1929, 8, 'progresse 'provoking': 959, 'punish' 0, 'quaid': 181, 'racial': g': 1993, 'ra 2004, 'realit ecently': 201 rences': 2024	peaking': 1815, 'pearls': 1816, 'peculiarity': 1817, 'pedestal': 1818, 'pencil': 1819, 'people': 1820, 'person': 1820, 'person': 1821, 'perfecte': 1822, 'perfecte': 1823, 'perfecte': 1825, 'performances': 1826, 'perhaps': 1827, 'period': 1828, 'perplexing': 1829, 'person': 1830, 'personalities': 1831, 'personally': 1832, 'peter': 1833 (antasm': 1835, 'phenomenal': 1836, 'philippa': 1837, 'phony': 1838, 'photography': 1839, 'photography': 1840, 'phrase': 1841, 'physical': 1842, 'pi': 1843, 'pice': 1845, 'pictures': 1846, 'picere': 1847, 'pieces': 1848, 'pile': 1849, 'pillow': 1850, 'pitch': 1851, 'pitiful': 1852, 'pixar': 1853, 'place': 1854, 'place': 1854, 'place': 1857, 'planned': 1858, 'place': 1859, 'playe': 1860, 'player': 1863, 'players': 1863, 'playing': 1864, 'plays': 1865, 'pleasant': 1866, 'pleasari': 1866, 'pleasari': 1867, 'poet': 1871, 'pier': 1872, 'plot': 1873, 'plug': 1874, 'plus': 1875, 'pm': 1876, 'poet': 1877, 'poetry': 1879, 'point': 1888, 'pointilistic': 1881, 'pointless': 1882, 'poised': 1883, 'poler': 1884, 'political': 1885, 'politically': 1886, 'politics': 1887, 'poorty': 1889, 'portrayal': 1890, 'popcorn': 1891, 'popcorn': 1891, 'power': 1992, 'portrayal': 1893, 'portrayals': 1894, 'portrayed': 1895, 'portraying': 1896, 'postive': 1897, 'possible': 1898, post': 1900, 'potted': 1991, 'power': 1992, 'previous': 1914, 'prepared': 1915, 'presence': 1916, 'presents': 1917, 'precisely': 1908, 'process': 1930, 'produce': 1931, 'produced': 1932, 'producer': 1933, 'profucers': 1934, 'product': 1935, 'production': 1936, 'professionals': 1937, 'professor' 1949, 'process': 1930, 'promote': 1940, 'prompted': 1941, 'prone': 1942, 'propaganda': 1943, 'properly': 1944, 'proud': 1945, 'producy': 1946, 'provided': 1947, 'provokes': 1949, 'ps': 1950, 'pseudo': 1951, 'psychological': 1952, 'psychotic': 1953, 'public': 1954, 'pull': 1955, 'pulling': 1967, 'puzzle': 1968, 'purce': 1968, 'pu
ing': 2043, 'nowned': 2053 e': 2063, 're t': 2073, 're 'rickman': 20 'road': 2095, 2106, 'room': s': 2118, 'ru 'salesman': 2 ave': 2141, ' ry': 2153, 's t': 2163, 'sc 'scripted': 2 84, 'secondar emi': 2196, ' ce': 2206, 's ral': 2217, ' 7, 'shattered 2238, 'shooti	n': 2034, 'relations': 2035, 'relationship': 2036, 'relationships': 2037, 'relatively': 2038, 'relaxing': 2039, 'release': 2040, 'released': 2041, 'relief': 2042, remaining': 2044, 'remake': 2045, 'remake': 2045, 'remake': 2045, 'remake': 2045, 'remake': 2056, 'repeated': 2056, 'repeated': 2056, 'repeating': 2057, 'repeats': 2058, 'repertory': 2059, 'reporter': 2060, 'represents': 2061, 'require': 2062, 'rescarched': 2064, 'resounding': 2065, 'respecting': 2066, 'rest': 2067, 'restrained': 2068, 'result': 2069, 'results': 2070, 'resume': 2071, 'retarded': 2072, 'restrained': 2068, 'revening': 2075, 'revenge': 2076, 'revere': 2077, 'reverse': 2078, 'reviewe': 2079, 'reviewer': 2080, 'reviewers': 2081, 'reviewes': 2082, 'rice': 2078, 'reviewe': 2089, 'rise': 2089, 'reviewer': 2080, 'reviewers': 2081, 'reviewes': 2082, 'rice': 2078, 'ridiculous': 2085, 'ridiculousness': 2086, 'right': 2087, 'riot': 2088, 'rips': 2089, 'rise': 2099, 'rita': 2091, 'rivalry': 2092, 'riveted': 2093, 'rize': 2070, 'rosevelt': 2096, 'robotic': 2097, 'rocked': 2099, 'rocked': 2099, 'rocks': 2100, 'roeg': 2101, 'role': 2102, 'roles': 2103, 'roller': 2104, 'rolls': 2105, 'roman 2107, 'roosevelt': 2108, 'roth': 2109, 'rough': 2110, 'round': 2111, 'routine': 2112, 'row': 2113, 'rpg': 2114, 'ropger': 2115, 'rubbish': 2116, 'rubin': 2117, 'in': 2119, 'running': 2120, 'ruthless': 2121, 'ryan': 2122, 'ryans': 2123, 'sabotages': 2124, 'sack': 2125, 'sacrifice': 2126, 'sad': 2127, 'said': 2128, 'sake': 2139, 'sawn': 2131, 'sample': 2132, 'sand': 2133, 'sandra': 2134, 'sappiest': 2135, 'sarcophage': 2136, 'sat': 2137, 'satanic': 2138, 'savalas': 2139, 'savant': 2150, 'scriene': 2154, 'scribling': 2155, 'schilling': 2155, 'schilling': 2156, 'scream': 2157, 'scneel': 2157, 'scame': 2148, 'scared': 2149, 'screenwriter': 2171, 'scrimm': 2172, 'script': 2174, 'scripting': 2175, 'scripts': 2176, 'screen': 2168, 'screend': 2169, 'screenplay': 2180, 'seener': 2181, 'seaen': 2181, 'seaen': 2183, 'secon': 2197, 'senese': 2198, 'senese': 2198, 'senesi': 2209,
1, 'site': 22 82, 'slightes iling': 2294, ying': 2305, 5, 'sorrentin 'space': 2327 2338, 'spiffy 'stage': 2350 arts': 2361, y': 2371, 'st s': 2382, 'st 'strident': 2 s': 2403, 'st s': 2414, 'su cess': 2424, days': 2435, 4, 'supposed1 454, 'surroun dney': 2465,	re': 2261, 'sing': 2262, 'singing': 2263, 'single': 2264, 'sinister': 2265, 'sink': 2266, 'sinking': 2266, 'sister': 2268, 'sisters': 2268, 'sisters': 2269, 'site': 2270, 'sitcoms': 2727, 'sites': 2273, 'sits': 2274, 'situation': 2275, 'situations': 2276, 'skilled': 2277, 'skip': 2278, 'slackers': 2279, 'slavic': 2280, 'slee': 2281, 'slightly': 2284, 'slimy': 2285, 'sloppy': 2286, 'slow': 2287, 'slurs': 2288, 'smack': 2289, 'small': 2299, 'smart': 2291, 'smells': 2292, 'smile': 2292, 'smith': 2295, 'smoothly': 2296, 'snider': 2297, 'snow': 2298, 'soap': 2299, 'sobering': 2300, 'social': 2301, 'soldiers': 2302, 'sole': 2303, 'solid': 2304, 'solid': 2304, 'solid': 2304, 'solid': 2306, 'sory': 2317, 'sort': 2318, 'soul': 2319, 'sound': 2329, 'sounded': 2321, 'sounds': 2322, 'soundtrack': 2323, 'sour': 2324, 'south': 2325, 'southern': ', 'spacek': 2328, 'spacey': 2329, 'span': 2330, 'speak': 2331, 'speaking': 2332, 'special': 2333, 'speed': 2334, 'spend': 2335, 'spent': 2336, 'spent': 2337, 'sph': '2339, 'splendid': 2340, 'spock': 2344, 'spoiled': 2343, 'spoiler': 2344, 'spoilers': 2345, 'spot': 2346, 'spy': 2347, 'squibs': 2348, 'stable': '0, 'stagy': 2351, 'stand': 2362, 'stamout': 2352, 'standout': 2353, 'stanwyck': 2354, 'star': 2355, 'starlet': 2356, 'starring': 2357, 'stars': 2358, 'start': 2359, 'started': 2362, 'stay': 2363, 'stay': 2363, 'stay': 2363, 'stay': 2363, 'stay': 2364, 'stayl': 2364, 'stayl': 2364, 'stayl': 2364, 'stayl': 2365, 'stayl': 2366, 'stayl': 2365, 'stayl': 2366, 'stayl': 2366, 'stayl': 2367, 'stayl': 2368, 'stayl': 2369, 'stroyetling': 2365, 'stayl': 2368, 'stroyetling': 2365, 'stayl': 2368, 'stroyetl': 2369, 'stroyetl': 2381, 'stockings': 2379, 'stroyetl': 2369, 'stro
e': 2498, 'te 'terrible': 2 519, 'themes' h': 2530, 'th t': 2540, 'th eless': 2551, 2, 'tolerable p': 2574, 'to s': 2585, 'to s': 2595, 'tr 'tremendous': k': 2616, 'tr 2627, 'turn': 2639, 'ue': 2 alled': 2650, 658, 'underra edly': 2667, y': 2676, 'un 5, 'unmatched	naches': 2488, 'team': 2489, 'tear': 2490, 'tears': 2491, 'technically': 2492, 'teddy': 2493, 'tedium': 2494, 'teen': 2495, 'teenagers': 2496, 'teerh': 2497, 'tel elevision': 2499, 'tell': 2500, 'terrific': 2511, 'terror': 2512, 'thi: 2513, 'thanks': 2514, 'theater': 2515, 'theatere': 2516, 'theatere': 2516, 'theatere': 2517, 'theatrical': 2518, 'the electric 2520, 'therapy': 2521, 'thick': 2522, 'thing': 2523, 'things': 2524, 'think': 2525, 'thinking': 2526, 'thomerson': 2527, 'thoroughly': 2528, 'thorsen': 2529, 'dought': 2531, 'thoughts': 2532, 'thousand': 2533, 'thread': 2534, 'three': 2535, 'threshold': 2536, 'threiled': 2537, 'thriller': 2538, 'thriller': 2538, 'throughly': 2542, 'thug': 2543, 'thumper': 2544, 'thumderbirds': 2545, 'thus': 2546, 'ticker': 2547, 'tickets': 2548, 'timiler': 2538, 'timer': 2559, 'timer': 2552, 'timers': 2552, 'timers': 2554, 'timers': 2554, 'timer': 2555, 'timer': 2554, 'timer': 2556, 'timer': 2556, 'timer': 2557, 'tolerate': 2561, 'tolerate': 2560, 'tong': 2566, 'tong': 2577, 'tolerate': 2577, 'total': 2578, 'total': 2578, 'total': 2578, 'total': 2578, 'total': 2578, 'total': 2578, 'total': 2579, 'tong': 2579
'uplifting': al': 2715, 'u 5, 'verbal': 'view': 2737, isual': 2748, yage': 2759, 0, 'warmth': 'watch': 2782 s': 2793, 'wb 4, 'went': 28 5, 'wholesome 2827, 'wise': 8, 'wonderful 2850, 'worry' 'wrap': 2862, 873, 'years': def transform '''This f	2703, 'upper': 2704, 'ups [*] : 2705, 'uptight': 2706, 'ursula': 2707, 'us': 2708, 'use': 2709, 'used': 2710, 'user': 2711, 'uses': 2712, 'using': 2713, 'usr': 2714, 'tter': 2716, 'utterly': 2717, 'valentine': 2718, 'value': 2719, 'values': 2720, 'vampire': 2721, 'vandiver': 2722, 'variation': 2723, 'vehicles': 2724, 'ventura' 2726, 'verbatim': 2727, 'versatile': 2728, 'version': 2729, 'versus': 2730, 'vessel': 2731, 'veteran': 2732, 'vey': 2733, 'vibe': 2734, 'victor': 2735, 'video': 'viewer': 2738, 'viewing': 2739, 'views': 2740, 'villain': 2741, 'villains': 2742, 'violence': 2743, 'violin': 2744, 'virtue': 2745, 'virus': 2746, 'vision': 27 'visually': 2749, 'vitally': 2759, 'vivian': 2751, 'vivid': 2752, 'vocal': 2753, 'voice': 2754, 'volatile': 2755, 'volcano': 2756, 'vomit': 2757, 'vomited': 275 'vulcan': 2760, 'waiting': 2761, 'waitress': 2762, 'walk': 2763, 'walked': 2764, 'wall': 2765, 'want': 2766, 'wanted': 2767, 'wanting': 2768, 'wants': 2769, 'wartine': 2771, 'warn': 2772, 'warning': 2773, 'wartime': 2774, 'warts': 2775, 'washed': 2776, 'washing': 2777, 'waste': 2778, 'wasted': 2779, 'waster': 2780, 'wating': 2784, 'watchable': 2783, 'watched': 2784, 'watching': 2785, 'watching': 2785, 'watchable': 2783, 'watched': 2784, 'watching': 2785, 'watchable': 2783, 'watched': 2784, 'watching': 2785, 'watchable': 2789, 'waster': 2796, 'wanting': 2791, 'wayne': 2792, 'vextina': 2809, 'walkever': 2806, 'watchable': 2807, 'whenever': 2808, 'whether': 2809, 'whine': 2810, 'whiny': 2811, 'white': 2812, 'whites': 2813, 'whoever': 2814, 'whole': 2816, 'wide': 2817, 'widmark': 2818, 'wife': 2819, 'wih': 2820, 'wild': 2821, 'wilkinson': 2822, 'william': 2823, 'willie': 2824, 'wily': 2825, 'win': 2826, 'win': 2826, 'worke': 2836, 'worke': 2846, 'worked': 2847, 'working': 2848, 'workes': 2849, 'w: 2851, 'worse': 2852, 'worst': 2863, 'worth': 2854, 'worthless': 2855, 'worthwhile': 2866, 'wrothe': 2866, 'worthy': 2867, 'would': 2888, 'wouldnt': 2859, 'woven': 2860, 'wow': 'write': 2863, 'writer': 2866, 'worth': 2866, 'writ
globa numbe #1:Ca freq= frequ for r f #In h frequ #2:Ca idf={ vocab for r	
#3:ge vocab vocab ind=0 for w #4:ca numbe csr=c #help	<pre>if len(word)<2: continue idf[word]=1+math.log((1+number_of_docs)/(1+frequency_of_words[word])) dict(sorted(idf.items(), key=lambda x:x[1], reverse=True)[:50]) #top idf 50 values etting the vocab for top 50 idf values 0={} 0={} 01={}</pre>
retur else: print : x=transform(c print(x) → Make sure t	co print instead of raw calling the function orm(corpus,vocab)
(2, 5) (2, 6) (4, 7) (5, 8) (7, 9) (7, 10) (9, 11) (9, 12) (9, 13) (10, 14) (10, 15) (11, 16) (12, 17) (15, 18) (16, 19) (17, 20) (17, 21) (18, 22) (19, 23) (19, 24)	0.5773502691896258 0.5773502691896258 1.0 1.0 0.7071067811865475 0.7071067811865475 0.5773502691896257 0.5773502691896257 0.5773502691896257 0.7071067811865476 0.7071067811865476 1.0 1.0 1.0 1.0 0.7071067811865475 0.7071067811865475 0.7071067811865475 0.7071067811865475 0.7071067811865475 0.7071067811865475 0.7071067811865475 1.0 0.1690308509457033 0.1690308509457033
(19, 25) (19, 26) (19, 27) (19, 28) (19, 30) (19, 31) (19, 32) (19, 33) (19, 34) (19, 35) (19, 36) (19, 37) (19, 38) (19, 39) (19, 40) (19, 41) (19, 42) (19, 43) (19, 44)	0.1690308509457033 0.1690308509457033
<pre>→ Make sure t e.g, x=transfo print(x[0] 100% [[0.57735027 0. 0. 0. 0. 0. 0.</pre>	se() nape:",y.shape) to print instead of raw calling the function orm(corpus,vocab)
0. 0. shape: (1, 50 print("vocab[for key,value print(f")	0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1]
14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35	baby 6.922918004572872 owls 6.922918004572872 florida 6.922918004572872 muppets 6.922918004572872 person 6.922918004572872 overdue 6.922918004572872 post 6.922918004572872 practically 6.922918004572872 practically 6.922918004572872 cross 6.922918004572872 helms 6.922918004572872 nerves 6.922918004572872 ladies 6.922918004572872 tited 6.922918004572872 lived 6.922918004572872 choices 6.922918004572872 dodge 6.922918004572872 blue 6.922918004572872 dodge 6.922918004572872 freeman 6.922918004572872