Count Based Distributional Semantics In this notebook we will implement a simple count based model. We will not try to optimize the model in order to get the best results possible. Instead we will try to keep things as simple as possible. Thus the main principles can become clear. We will use a mid-sized corpus as a compromise between fast processing and usefull results. In this notebook we will use German data. However, not much knowledge of German is required. And learning some new German words on the fly is not too bad. Caution Some of the cells require quite long computation time! Data We use a corpus of 300k sentences from wikipedia collected by the university of Leipzig. You can download the required data from http://wortschatz.uni-leipzig.de/en/download/ . The file used here is deu_wikipedia_2016_300K-sentences Reading the data The texts in the corpus are already split into sentences. We read thes sentence and word-tokenize each sentence. To save time we use the infected words and do not do anly lemmatization or stemming. In [9]: import codecs import nltk sentences = [] #Caution: change the path to this file! source = codecs.open('deu news 2016 300K/deu news 2016 300K-sentences.txt','r') for line in source: nr,sent = line.split('\t') sentences.append(nltk.word tokenize(sent.strip(),language='german')) Let us check whether the senteces are in the list as we expect them to be: In [10]: sentences[447] ['Aber', Out[10]: 'auch', 'der', 'erst', 'im', 'Juli', 'eingeweihte', 'Deli', 'Market', ',', 'der', 'Shopping', 'del', 'Sol', 'finden', 'ist', 'und', 'hochwertige', 'Lebensmittel', 'anbietet', ',', 'hatte', 'mit', 'einem', 'defekten', 'Dach', 'zu', 'kämpfen', 'welches', 'teilweise', 'einbrach', 'und', 'ebenfalls', 'für', 'interne', 'Überschwemmungen', 'sorgte', Vectors of Co-occurrence Values We write a function that computes vectors with co-occurrence numbers for a number of words with a given list of 'context' words. As context words we take all words that exceed a specified minimum frequency. KCo-occurrence with rare words will hardly contribute to the comparison between context vectors. Another parameter that needs to be set is the maximum distance between two words to count as co-competition, known as the window size. We proceed sentence by sentence here. This means that the last word of a sentence and the first word of the next sentence are never counted as co-occurring. Note, however, that sentence bondaries are often not taken into account in such procedures. It is not clear whether this has a seriosu impact. Another detail to note is that we first calculate the window size for the competition, and then determine the relevant words in the window. Alternatively, you can first remove all irrelevant words and then determine the words within the window. In general, a larger window can compensate for a too small corpus. In addition, a smaller window captures more syntactical properties of a word, while a larger window takes into account broader semantic relationships. When all co-occurrence values are calculated, we normalize the length of the vectors. In [11]: from scipy import sparse from collections import Counter import numpy as np def mag(x): return np.sqrt(x.dot(x)) def makeCV(words, sentences, window = 2, minfreq = 10): #first count all words freq = Counter() for s in sentences: freq.update(s) #determine which words have to be used as features or context words context words = [w for w, f in freq.items() if f > minfreq] #we add all words to the context words, if they are not already in that list. #Just for convenience to have all words in one list. for w in words: if w not in context words: context words.append(w) dim = len(context words) #Give each context word a unique number and build dictionaries to switch between numbers and words $cw2nr = {}$ for i in range(dim): cw = context words[i] cw2nr[cw] = i $w2nr = \{\}$ for i in range(len(words)): w = words[i] w2nr[w] = i#initialize a matrix #We use a sparse matrix class from scipy matrix = sparse.lil_matrix((len(words),dim)) #Now we start the real work. We iterate through all sentences and count the co-occurrences! for s in sentences: i s = [cw2nr.qet(w,-1) for w in s]for i in range(len(s)): w = s[i]if w in words: i w = w2nr[w]for j in range(max(0,i-window),min(i+window+1,len(s))): if i != j: # a word is not in its own context! i cw = i s[j]**if** i cw > 0: matrix[i w,i cw] += 1 #finally make a dictionary with vectors for each word wordvectors = {} for w,i w in w2nr.items(): v = matrix[i w].toarray()[0] wordvectors[w] = v/mag(v) return wordvectors Let us test the function for a short list of words: In [19]: vectors_count = makeCV(['Kloster','Kirche','Garten','Haus','Hof','Schweden','Deutschland','betrachten','anscha array([0., 0., 0., ..., 0., 0., 0.]) Out[19]: Since the vectors have the same length, we can use the inner product, which is identical to the cosine, for comparison: In [13]: print(vectors_count['Schweden'].dot(vectors count['Deutschland'])) print(vectors_count['Schweden'].dot(vectors_count['Hof'])) 0.820131020955634 0.6369476065012749 Next we want to know, what words are most similar to a given word. To do so, we need to compare a word with each other word in the list. We use an ordered list to store the results. Since these list always sort ascending, we need to consider always the last elements of this list. Finally, we return the results in inverse order. In [14]: import bisect def most similar(word, vectors, n): best = [] vec w = vectors[word] for z in vectors: if z == word: continue sim = vec w.dot(vectors[z]) #we have to add this result only, if we do not yet have n results, or if the similarity is larger than if len(best) < n or sim > best[0][0]: bisect.insort(best, (sim, z)) best = best[-n:] return best[::-1] #present the list in descending order In [15]: most similar('Garten', vectors count, 3) Out[15]: [(0.7600337758643444, 'Hof'), (0.7175335523650366, 'Haus'), (0.6911270544850551, 'Kloster')] It is more interesing to find the most similar word if we have more words to choose from. Let us collect some mid-frequency words to do so. In [16]: wortfrequenz = Counter() for satz in sentences: wortfrequenz.update(satz) vocabulary = [w for w,f in wortfrequenz.items() if 30 < f < 3000]</pre> vocabulary size = len(vocabulary) print(vocabulary size) 11217 Now it takes some time to compute all vectors. In [42]: vectors_count = makeCV(vocabulary, sentences, window=2) In [20]: most_similar('Garten', vectors_count, 10) Out[20]: [(0.7600337758643444, 'Hof'), (0.7175335523650366, 'Haus'), (0.6911270544850551, 'Kloster'), (0.594242060096992, 'Schweden'), (0.5843819971575435, 'anschauen'), (0.5241629791086438, 'betrachten'), (0.5155051637362442, 'Deutschland'), (0.5113054285912206, 'beobachten'), (0.4682522543097664, 'Kirche')] In [21]: most_similar('betrachten', vectors_count, 10) Out[21]: [(0.9283079870321336, 'beobachten'), (0.7048240848026528, 'anschauen'), (0.5296872135718191, 'Schweden'), (0.5241629791086438, 'Garten'), (0.4893808439018669, 'Haus'), (0.4761698863559131, 'Kirche'), (0.4346370889054621, 'Hof'), (0.4113083216746324, 'Deutschland'), (0.36126647423022834, 'Kloster')] In [22]: most_similar('Schweden', vectors count, 10) Out[22]: [(0.820131020955634, 'Deutschland'), (0.6698001763000421, 'Kirche'), (0.6445428851191304, 'Haus'), (0.6396984520026665, 'anschauen'), (0.6369476065012749, 'Hof'), (0.594242060096992, 'Garten'), (0.5296872135718191, 'betrachten'), (0.517356880541151, 'beobachten'), (0.4333519064039629, 'Kloster')] **Evaluation** The examples look nice, but how good are our vectors? There are a number of tests that can be done to check the quality of the vectors. One type of test is to compare the calculated similarity between two words to the similarity that subjects have given for word pairs. We can simply use the correlation to check the similarity. A dataset with similarity assessments for German word pairs is the so-called Gur350 dataset, which was developed by the working group of Prof. Dr. Iryna Gurevych at the TU Darmstadt. Further information and a link for the download can be found here: https://www.informatik.tu-darmstadt.de/ukp/research_6/data/semantic_relatedness/german_relatedness_datasets/index.en.jsp The data does not give a similarity between the words but a relationship, which is not exactly the same. For the evaluation, we only use the word pairs for which both words are contained in our vocabulary. In [25]: import math testfile = codecs.open('../../TMNB/Corpora/wortpaare350.gold.pos.txt','r','utf8') testfile.readline() testdata = [] missing = set()for line in testfile: w1, w2, sim, p1, p2 = line.split(':')if w1 in vocabulary and w2 in vocabulary: testdata.append((w1,w2,float(sim))) def evaluate(data, vectors): gold = []predicted = [] for v,w,sim in data: pred = vectors[v].dot(vectors[w]) gold.append(sim) predicted.append(pred) $\#print(v, w, pred, sim, sep = ' \ t')$ av p = sum(predicted)/len(predicted) $av_g = sum(gold)/len(gold)$ cov = 0var g = 0var p = 0for s,t in zip(gold,predicted): $cov += (s-av_g) * (t-av_p)$ $var_g += (s-av_g) * (s-av_g)$ $var_p += (t-av_p) * (t-av_p)$ return cov / math.sqrt(var g*var p) FileNotFoundError Traceback (most recent call last) Input In [25], in <cell line: 2>() 1 import math ---> 2 testfile = codecs.open('../../TMNB/Corpora/wortpaare350.gold.pos.txt','r','utf8') 3 testfile.readline() 5 testdata = [] File /opt/anaconda3/lib/python3.8/codecs.py:905, in open(filename, mode, encoding, errors, buffering) 901 if encoding is not None and \ 902 'b' not in mode: # Force opening of the file in binary mode
mode = mode + 'b' --> 905 file = builtins.open(filename, mode, buffering) 906 if encoding is None: return file FileNotFoundError: [Errno 2] No such file or directory: '../../TMNB/Corpora/wortpaare350.gold.pos.txt' In [18]: len(testdata) 119 Out[18]: Note, that we are using only a very small part of the test data from the Gur350 dataset for testing. Thus we cannot compare our results to official results for these data! We can of course include more words in our vocabulary of mid-frequency words, but many words simpy are not in our corpus. In [19]: evaluate(testdata, vectors count) 0.17241440456996857 Out[19]: Despite the good looking lists that we generated above, we see that the calculated similarities correlate only very weakly with the similarity judgments. Somewhat more advanced Our first attempt still offers various possibilities for optimization. First of all we could not use the simple co-occurrence frequencies but the Positive Pointwise Mutual Information (PPMI). The Pointwise Mutual Information between two words a and b is in principle the ratio between the actual probability that they will occur together and the expected probability that they will occur if their occurrences were independent. We define $pmi(a,b) = log(\frac{p(ab)}{p(a)p(b)})$. The ppmi(a, b) is now the pmi(a, b) if this value is positive, and otherwise 0. The use of the ppmi is based on the assumption that only the fact that words occur together more often than expected is interesting, while lower probabilities are more likely to be based on coincidence or in any case do not say anything interesting about word pairs. Another problem with our naive approach was that the vectors are extremely long. The vectors not only need a lot of storage space. Many of the dimensions also contain little or only redundant information. This problem can be solved by using a common dimension reduction process. Below we will use Singular Value Decomposition for this and then use the most important 100 dimensions. In [46]: **from** sklearn.decomposition **import** TruncatedSVD import math import numpy as np def makeCV SVD(words, sentences, window = 2, minfreq = 10, size = 256): freq = Counter() for s in sentences: freq.update(s) context_words = [w for w,f in freq.items() if f > minfreq] for w in words: if w not in context words: context words.append(w) dim = len(context words) $cw2nr = {}$ for i in range(dim): cw = context words[i] cw2nr[cw] = i $w2nr = {}$ for i in range(len(words)): w = words[i] w2nr[w] = imatrix = sparse.lil_matrix((len(words),dim)) n = 0for s in sentences: n+=1i s = [cw2nr.get(w,-1) for w in s]for i in range(len(s)): w = s[i]if w in words: for j in range(max(0,i-window),min(i+window+1,len(s))): **if** i != j: $i_cw = i_s[j]$ **if** i cw > 0: $matrix[i_w,i_cw] += 1$ #up to here nothing new N = matrix.sum() #Let us get the probabilities for each word: freq w = matrix.sum(axis = 1)freq_w = np.array(freq_w.T)[0] prob_w = np.array(freq_w) / N #Let us get the probabilities for each context word: freq cw = matrix.sum(axis = 0)freq_cw = np.array(freq_cw)[0] prob_cw = np.array(freq_cw) / N (rows, cols) = matrix.nonzero() #Returns a tuple of arrays (row, col) containing the indices of the non-zero for i_w,i_cw in zip(rows,cols): p = matrix[i w, i cw]/N $p_w = prob_w[i_w]$ p_cw = prob_cw[i_cw] $ppmi = max(0, math.log(p/(p_w * p_cw)))$ matrix[i_w,i_cw] = ppmi #We cannot be completely sure, that for every word we found at least some positive pmi-values. #The frequent neighbors of a word eventaully are not in the set of context words #A row with only zeros, will cause a problem for the SVD for i in range(len(words)): if np.sum(matrix[i]) print("Empty row for:", words[i]) print("Please remove this word from the wordlist.") svd = TruncatedSVD(n_components=size) svd.fit(matrix) matrix = svd.transform(matrix) wordvectors = {} for w,i_w in w2nr.items(): $v = matrix[i_w]$ wordvectors[w] = v/mag(v)return wordvectors In [47]: vectors_SVD = makeCV SVD(vocabulary, sentences, window=2, size=100) In [48]: most_similar('Garten', vectors SVD, 10) [(0.7810250655750048, 'Körper'), Out[48]: (0.7793117054291634, 'Laden'), (0.7593780649332582, 'Hund'), (0.7536957161666307, 'Geschmack'), (0.7529221520215187, 'Beruf'), (0.7332369249730386, 'Balkon'), (0.7317826842138249, 'Charme'), (0.7295553944972246, 'Mund'), (0.7276913223621029, 'Bett'), (0.7186351640200018, 'Hof')] In [49]: most_similar('betrachten', vectors SVD, 10) Out[49]: [(0.810423435388816, 'agieren'), (0.7877999287255073, 'orientieren'), (0.7791489152019035, 'bewegen'), (0.7608466529178798, 'aufhalten'), (0.7576881952580662, 'ignorieren'), (0.7553047730854288, 'bezeichnen'), (0.7550711017117226, 'organisieren'), (0.7546950738136995, 'gefährden'), (0.7479776257602084, 'vorgehen'), (0.74724751977293, 'beeinflussen')] In [50]: most_similar('Schweden', vectors SVD, 10) [(0.8976080862395008, 'Spanien'), Out[50]: (0.8937456150121426, 'Dänemark'), (0.8914081101850169, 'Norwegen'), (0.8886603831524063, 'Tschechien'), (0.8846036734853973, 'Finnland'), (0.8769569275332175, 'Ungarn'), (0.8712547733396998, 'Portugal'), (0.8707173048342136, 'Belgien'), (0.8692788751387291, 'Italien'), (0.8633694772872089, 'Irland')] In [51]: most_similar('Park', vectors_SVD, 10) [(0.7840628615147447, 'Tower'), Out[51]: (0.7578589097604853, 'State'), (0.7233800759562373, 'Club'), (0.7229083946474175, 'National'), (0.7201041868251652, 'School'), (0.7076997520310662, 'Chicago'), (0.7031226315548857, 'South'), (0.701415495916093, 'Hotel'), (0.6989094983260197, 'New'), (0.697688501084682, 'Air')] We notice that park is not interpreted as a synonym for garden, but rather is perceived as part of English place names. The context 'English words' has apparently become more important. In [52]: evaluate(testdata, vectors SVD) Traceback (most recent call last) Input In [52], in <cell line: 1>() ----> 1 evaluate (testdata, vectors_SVD) NameError: name 'evaluate' is not defined The optimization was apparently successful: the correlation has become significantly larger. Bullinaria and Levy (2007), Mohammad and Hirst (2012), Bullinaria and Levy (2012) and Kiela and Clark (2014) provide overviews of the combinations of parameters, training quantities, etc. that lead to optimal results. **Exercise:** Reading the STW German Synonyms dataset import numpy as np In [65]: from sklearn import metrics stw = codecs.open ('STW-Syn.txt', 'r', 'utf8') data stw = []for line in stw : #Splitting to get the words and their label $w1 , w2 , label = line.strip().split('\t')$ #if words are in vocabulary then append if w1 in vocabulary and w2 in vocabulary : data_stw.append ((w1 , w2 , label)) In [66]: #Checking the similarity rankings for different words most_similar('Global', vectors_SVD, 10) [(0.8731351048734526, 'Research'), Out[66]: (0.8685843804181166, 'Group'), (0.8605837909793592, 'Energy'), (0.8515605888317561, 'Capital'), (0.8454786942933771, 'Inc'), (0.8406273546449221, 'Europe'), (0.8405246088673665, 'Technology'), (0.8363250304183301, 'Digital'), (0.8178653749201173, 'President'), (0.8165933185889089, 'Network')] In [67]: most_similar('Landkreis', vectors_SVD, 10) [(0.8024150082814971, 'Kreis'), Out[67]: (0.7406383629646376, 'Lörrach'), (0.7392727828649166, 'Aurich'), (0.7122872947181331, 'Leer'), (0.6976964299241158, 'Bezirk'), (0.693245271463566, 'Gießen'), (0.6843357747300666, 'Weiden'), (0.6842304638668038, 'Emden'), (0.6799150808902322, 'Böblingen'), (0.6788010942358316, 'Dorf')] In [68]: most_similar('Hund', vectors_SVD, 10) [(0.8037163385755053, 'Tier'), Out[68]: (0.7978506315523081, 'Rucksack'), (0.7713363976013471, 'Laden'), (0.7639887963472162, 'Junge'), (0.763077042706787, 'Baby'), (0.7605940760609347, 'Bett'), (0.7590656536489175, 'Hut'), (0.7357053275546623, 'Kind'), (0.7280352608354306, 'Freund'), (0.7279683834049742, 'tot')] **Evaluate the ranking using AUC** In [69]: #Function to evaluate using AUC def evaluate_auc(data, vectors): labels = []predictions = #for each element in stw_dataset for v , w , label in data : #Find the similarity value pred = vectors[v].dot(vectors[w]) labels.append(label) predictions.append(pred) #roc curve returns the false positive and true positive rates. #These are used to calculate AUC fpr , tpr , thresholds = metrics.roc_curve(labels , predictions , pos_label = '+') return metrics.auc(fpr , tpr) In [70]: evaluate_auc (data_stw , vectors_SVD) 0.9692460317460319 Out[70]: We see that the value is higher than the previous value. Varying the different attributes In [73]: #Changing the number of dimensions and window size vectors SVD = makeCV SVD(vocabulary, sentences, window=5, size=50) evaluate_auc (data_stw , vectors_SVD) 0.9598214285714285 Out[73]: In [74]: #Changing window size to 5 vectors SVD = makeCV SVD(vocabulary, sentences, window=5, size=100) evaluate_auc (data_stw , vectors SVD) 0.9692460317460319 Out[74]: In []: #Changing corpus size size = len(sentences) vocab_new = [w for w,f in wortfrequenz.items() if 15 < f < 2000]</pre> vectors SVD = makeCV SVD(vocab new,sentences[0:size-1000],window=2, size=100) for line in stw : #Splitting to get the words and their label $w1 , w2 , label = line.strip().split('\t')$ #if words are in vocabulary then append if w1 in vocab_new and w2 in vocab_new : data_stw.append ((w1 , w2 , label)) evaluate_auc (data_stw , vectors_SVD) In [77]: | #Reducing the vocabulary size vocabulary_new = [w for w,f in wortfrequenz.items() if 50 < f < 2500]</pre> vocabulary_len = len(vocabulary_new) print("Length of vocabulary-", vocabulary_len) vectors_SVD = makeCV_SVD(vocabulary_new,sentences,window=2, size=100) for line in stw : #Splitting to get the words and their label $w1 , w2 , label = line.strip().split('\t')$ #if words are in vocabulary then append if w1 in vocabulary_new and w2 in vocabulary_new : data_stw.append ((w1 , w2 , label)) evaluate_auc (data_stw , vectors_SVD) Length of vocabulary- 7343 Traceback (most recent call last) Input In [77], in <cell line: 16>() if w1 in vocabulary new and w2 in vocabulary new : data stw.append ((w1 , w2 , label)) Input In [69], in evaluate auc(data, vectors) 6 #for each element in stw dataset 7 for v , w , label in data : #Find the similarity value pred = vectors[v].dot(vectors[w]) labels.append(label) 10 11 predictions.append(pred) KeyError: 'Herrschaft' We see that the auc value changes depending on the varying attributes Reference: http://textmining.wp.hs-hannover.de/DistributionelleSemantik.html Literature Rubenstein, H., & Goodenough, J. B. (1965). Contextual Correlates of Synonymy. Commun. ACM, 8(10), 627-633. Ruge, G., & Schwarz, C. (1990). Linguistically based term associations—A new semantic component for a hyperterm system. Proceedings of 1st International ISKO-Confercence, Darmstadt, 88-96. Crouch, C. J. (1990). An approach to the automatic construction of global thesauri. Information Processing & Management, 5, 629– 640. Crouch, C. J., & Yang, B. (1992). Experiments in automatic statistical thesaurus construction. In Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval (pp. 77-88). ACM. Grefenstette, G. (1992). Use of syntactic context to produce term association lists for text retrieval. 89–97. Karlgren, J., & Sahlgren, M. (2001). From Words to Understanding. In Foundations of Real-World Intelligence (S. 294–308). CSLI Publications. Bullinaria, J. A., & Levy, J. P. (2007). Extracting semantic representations from word co-occurrence statistics: A computational study. Behavior research methods, 39(3), 510-526. Mohammad, S. M., & Hirst, G. (2012). Distributional measures of semantic distance: A survey. arXiv preprint arXiv:1203.1858. Bullinaria, J. A., & Levy, J. P. (2012). Extracting Semantic Representations from Word Co-occurrence Statistics: Stop-lists, Stemming and SVD. Behaviour Research Methods, 44(3), 890-907. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems (pp. 3111-3119). Kiela, D., & Clark, S. (2014). A Systematic Study of Semantic Vector Space Model Parameters. In 2nd Workshop on Continuous Vector Space Models and their Compositionality (CVSC). Levy, O., Goldberg, Y., & Dagan, I. (2015). Improving distributional similarity with lessons learned from word embeddings. Transactions of the Association for Computational Linguistics, 3, 211-225.