Implement a modified version of COALS algorithm

Implement a modified version of COALS (refer to this paper for COALS implementation) using the COVID corpus.

(a) Instead of the using correlation, use the ratio of the probabilities $\frac{p_{ik}}{p_{ik}}$ to obtain the co-

occurrence matrix

- (b) Try to limit the vocabulary size to ~7K to compute the word embedding.
- (c) Display the vocabulary size and the size of the matrix.
- (d) Identify five nouns and verbs (relevant to COVID19) from the corpus.
- (e) Generate five similar words for each (chosen in step b) and display the result with cosine distance for each one.
- (f) Use multi-dimensional scaling to visualize certain concepts (relevant to COVID 19) as shown in Fig.1. Take three concepts that you feel are relevant to this corpus and use a maximum of 10 words per concept. Make sure that the plot is clear and not cluttered

	NDIOT				
					HARM HA
	SHOULDER				BRAZIL
• E	ANANDE				∎ Moscow
• T	OE • LEG				
• TO	OTHOELD				■ HAWAII
	• N•/AEAD				∎токуо
○ MC	DUSE				- ATLANTA
	O DOG				
	0 OAT	o TUF	TIF		
		0.001	○ LION		■ NASHVILLE
	S KITTEN	000W			
				○OYSTER	
				Eig 1	
	○ BULL			Fig.1	

Note: Follow all the instructions given in Assignment 1.

You may submit these two assignments as one single assignment

	Tasks	Total marks
Assignment 2	Forming the cooccurrence matrix (a)	25 marks
Assignment 3	Similarity measure	5 marks
Assignment 3	Listing of similar words	5 (1 for each)
Assignment 3	Graph with clear visualisation	15 marks

Additional Information

Large matrices are very common in machine learning and NLP applications. There are several ways to handle large matrix. If the matrix is sparse, then you may try https://docs.scipy.org/doc/scipy/reference/sparse.html. Otherwise, you may use pandas or use h5py.

Same code is given for your convenience -You may use it, but at your own risk 😄

```
#To reduce the matrix size, you may want to reduce the
# vocabulary size.
freq_dist = FreqDist(tokens)
vocab = []
for word in tokens:
    if not stop_words(word) and freq_dist[word] > 20:
        vocab.append(word)
```