

MCA Assignment-3
Arushi Chauhan
2016019

Q1: Terminology:

1. Context and center word - it refers to the window consisting of n words to the right and left of a chosen center word. For instance, in the sentence "a lion crossed the road", if "crossed" is a center word and window size is 2, then context would contain the words ["a", "lion", "the", "road"].
2. Bag of words model - It takes a 'multiset' of words in a sentence, discarding the positional information and grammar.

Library used: Pytorch

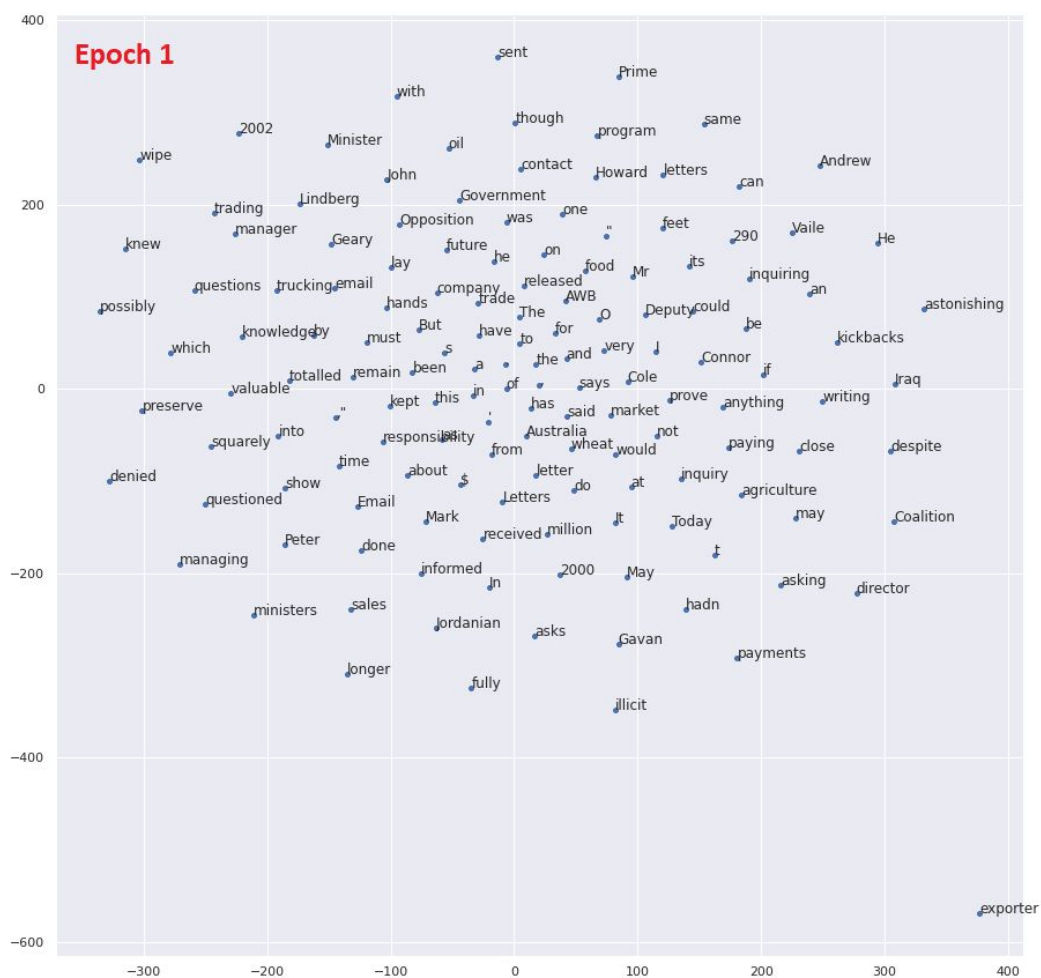
Word2Vec model:

The Word2Vec model consists of two components - CBOW (continuous bag of words) and skipgram. The CBOW part takes in context as input and predicts the center word. The order of the words in the context does not matter, that is why it is a bag of words model. The skipgram takes the center word as input and outputs the most relevant context of the words. Both of these are shallow models, i.e. they generally have one hidden layer to compute the vector representations.

Implementation details:

For this assignment, I have used one embedding layer and one hidden layer in both CBOW and SkipGram models. The models were trained for 50 epochs each. In case of CBOW model, the context words were given as inputs. These context words were passed through embedding layer and the mean of these embedded vectors was taken to get a mean representation of the context. This mean vector was then passed through an FC layer to get a vector of size the length of vocabulary. Taking log softmax of this vector gives the word which the CBOW model predicts as the center word as in one-hot encoding. For skipgram model, during training the center word was passed through the embedding layer and the rest of the training proceeded similar as above.

The TSNE plots for epochs 10, 30 and 50 are as shown at the end of report in figures 1,2, and 3. We can also see the compiled TSNE plots for each epoch in the following GIF. From these plots, we can notice that words that occur frequently together in the corpus move closer as epochs increase. For instance, in epoch 50 (figure 3), we can see ["Howard", "Program"], ["minister", "denied", "email"], ["possibly", "manager", "trading"] are occurring together.



Q2: The following values were reported after relevance feedback and query expansion were implemented:

MAP values	Relevance Feedback	Relevance Feedback with Query Expansion
1 iteration	0.6544	0.6864
3 iterations	0.9301	0.9720

Steps followed in relevance feedback:

1. Load the ground truth from the file MED.REL.
2. For each query, follow these steps:
 - a. Using similarity matrix passed as argument, compute the top 10 documents that are perceived as relevant (according to similarity computed of queries and documents)
 - b. Classify the perceived relevant documents as relevant and non-relevant using ground truth as comparison.
 - c. Sum all the relevant documents and non-relevant documents separately and multiply the sum with alpha and beta respectively.
 - d. Divide the above vectors with length of relevant and non-relevant documents respectively.
 - e. Update the query by adding the modified sum of relevant documents and subtracting the modified sum of non-relevant documents from it.

Steps followed in relevance feedback with query expansion:

1. Load the ground truth from the file MED.REL.
2. For each query, follow these steps:
 - a. From the ground truth, collect top n terms (those having highest TF-IDF values)
 - b. Add the TF-IDF values of these top n terms to the query.

Role of alpha and beta:

Alpha and beta control the sum of relevant and non-relevant documents that are added to the query. The following values were obtained for different values of alpha and beta:

Alpha	Beta	MAP Value
0.75	0.15	0.65
0.75	0	0.62
0.5	0.15	0.63
0	0.15	0.51
0.9	0.8	0.66

We can see that the influence of alpha is greater than the influence of beta on the MAP value, demonstrated by the huge gap in MAP values when alpha = 0 versus when beta = 0. Furthermore, the values tend to remain mostly near 0.64 for most values of alpha and beta pairs.

Q1 figures:

Figure 1

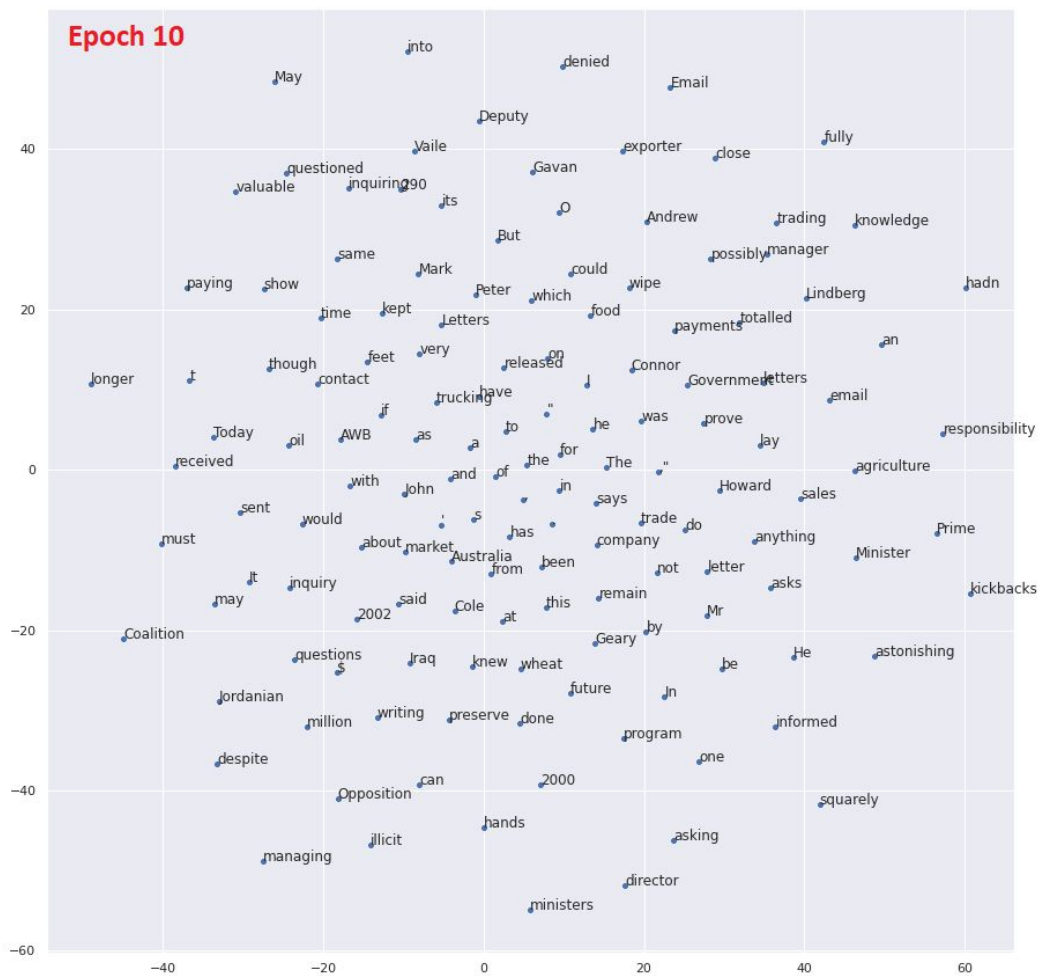


Figure 2



Figure 3

