

50.040 Natural Language Processing, Fall 2024

Homework 1

Due 04 October 2024, 23:59pm

Homework 1 will be graded by Chen Huang

Overview

Word embeddings are dense vectors that represent words, and are capable of capturing semantic and syntactic similarity, relation with other words, etc. We have introduced two approaches in the class to learn word embeddings: **Count-based** and **Prediction-based**. Here we will explore both approaches. Note that we use “word embeddings” and “word vectors” interchangeably.

Before we start, you need to download the text8 dataset. Unzip the file and then put it under the “data” folder. The text8 dataset consists of one single line of long text. Please do not change the data unless you are requested to do so.

Requirements

Python 3.5 or above / gensim / sklearn / numpy

1 Count-Based Word Embeddings

Co-Occurrence A co-occurrence matrix counts how often things co-occur in some environment. Given some word w_i occurring in the document, we consider the context window surrounding w_i . Supposing our fixed window size is n , then this is the n preceding and n subsequent words in that document, i.e. words w_{i-n}, \dots, w_{i-1} and w_{i+1}, \dots, w_{i+n} . We build a co-occurrence matrix M , which is a symmetric word-by-word matrix in which m_{ij} is the number of times w_j appears inside w_i 's window.

Example: Co-occurrence with fixed window of $n = 1$:

Document 1: “learn and live”

Document 2: “learn not and know not”

	and	know	learn	live	not
and	0	1	1	1	1
know	1	0	0	0	1
learn	1	0	0	0	1
live	1	0	0	0	0
not	1	1	1	0	0

Question 1.1 [written] (2 points) To have a better understanding of the co-occurrence matrix, please write a matrix M with fixed window of $n = 2$: Document 1: “learn to know” Document 2: “not learn to not know well”

Positive Pointwise Mutual Information (PPMI) Pointwise mutual information (PMI) is one of the most important concepts in NLP. The pointwise mutual information between a target word w and a context word c is defined as:

$$\text{PMI}(w, c) = \log_2 \left(\frac{P(w, c)}{P(w)P(c)} \right)$$

It is more common to use positive PMI (PPMI) which replaces all negative PMI values with zero. Given co-occurrence matrix $M \in \mathbb{Z}^{N \times N}$ of N words, m_{ij} is the element of i th row and j th column. The PPMI matrix can be calculated as:

$$\text{PPMI}_{ij} = \max \left(\log_2 \frac{p_{ij}}{p_{i*}p_{*j}}, 0 \right)$$

where

$$p_{ij} = \frac{m_{ij}}{\sum_{i=1}^N \sum_{j=1}^N m_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^N m_{ij}}{\sum_{i=1}^N \sum_{j=1}^N m_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^N m_{ij}}{\sum_{i=1}^N \sum_{j=1}^N m_{ij}}$$

For the details of PMI and PPMI, please refer to this link.

Principal Components Analysis (PCA) and Truncated Singular Value Decomposition (Truncated SVD) The rows (or columns) of co-occurrence matrix or PPMI matrix can be utilized as word vectors, but the vectors will be large in general (linear in the number of distinct words in a corpus). Thus, our next step is to run dimensionality reduction. In particular, we will first run PCA (Principal Components Analysis) to reduce the dimension. In practice, it is challenging to apply PCA to large corpora because of the memory needed to perform PCA. However, if you only want the top k vector components for relatively small k - known as Truncated SVD - then there are reasonably scalable techniques to compute those iteratively.

Question 1.2 [code] (2 points) Implement the function *distinct_words* that reads in *corpus* and returns distinct words that appeared in the corpus and the number of distinct words. Then, run the sanity check cell to check your implementation.

Question 1.3 [code] (6 points) Implement *compute_word_matrix* that reads in *corpus* and *window_size*, and returns a co-occurrence matrix, PPMI matrix and a word-to-index dictionary. Then, run the sanity check cell to check your implementation.

Question 1.4 [code] (5 points) Implement *dimension_reduction* function with python package sklearn. decomposition. Then, run the sanity check cell to check your implementation.

Question 1.5 [code] (5 points) Implement *plot_embeddings* function to visualize the word embeddings on a 2-D plane. Then, run the sanity check cell to check your implementation.

2 Prediction-Based Word Embeddings

Word2vec Word2vec is a software package that contains two algorithms named CBOW and skip-gram (Mikolov 2013). In the CBOW architecture, the model predicts the current word from a window of surrounding context words. In the continuous skip-gram architecture, the model uses the current word to predict the surrounding window of context words. The architectures are shown as Fig 1:

Question 2.1 [code] (3 points) Complete the code in the function *create_word_batch*, which can be used to divide a single sequence of words into batches of words. For example, the word sequence ["I", "like", "NLP", "So", "does", "he"] can be divided into two batches, ["I", "like", "NLP"], ["So", "does", "he"], each with batch size = 3 words. It is more efficient to train word embeddings on batches of word sequences rather than on a long single sequence. Then, run the sanity check cell to check your implementation.

Question 2.2 [code] (3 points) Use *Word2Vec* function to build a word2vec model. Please use the parameters we have set for you. It may take a few minutes to train the model.

Question 2.3 [code] (6 points) Implement *get_word2Ind* function first. Then, run the sanity check cell to check your implementation. Use *get_word2Ind*, *dimension_reduction*, and *plot_embeddings* functions to visualize the word embeddings of the first 300 words in the vocabulary.

Question 2.4 [code] (4 points) (1) Find the most similar words for the given words "dog", "car", "man". You need to use *model.wv.most_similar* function. (2) Find out which word will it be for x in the pairs woman : king :: man : x ? You need to use *model.wv.most_similar* function.

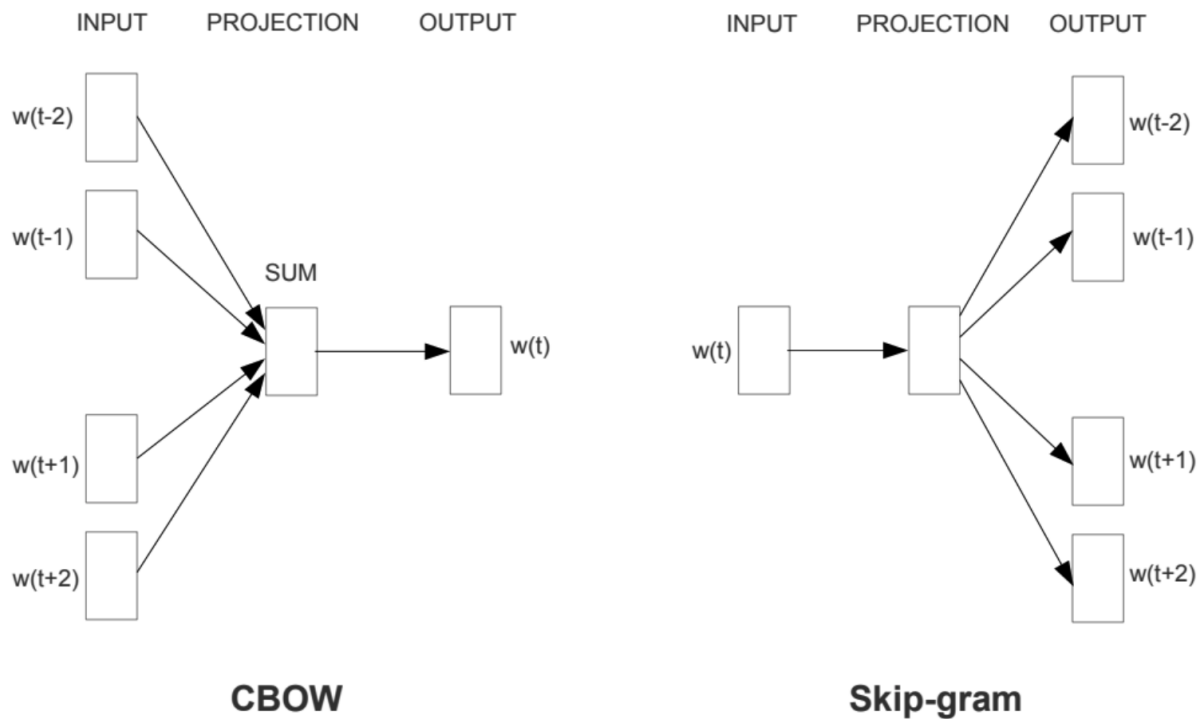


Figure 1: Word2vec Structure

Question 2.5 [code+written] (4 points) It's important to be cognizant of the biases (gender, race, sexual orientation etc.) implicit in our word embeddings. Bias can be dangerous because it can reinforce stereotypes through applications that employ these models. Use the *most_similar* function to find two cases where some bias is exhibited by the vectors. Please briefly explain the example of bias that you discover.

How to submit

1. Fill up your student ID and name in the Jupyter Notebook.
2. Click the Save button at the top of the Jupyter Notebook.
3. Select Cell - All Output - Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
4. Select Cell Run All. This will run all the cells in order, and will take several minutes.
5. Once you've rerun everything, select File – Download as – PDF via LaTeX.
6. Look at the PDF file and make sure all your solutions are there, displayed correctly. **The PDF is the only thing your graders will see!**
7. Submit your PDF on eDimension.