

# Learning Sense Embeddings

Michał Ostyk-Narbutt (1854051)

Prof. Roberto Navigli

Natural Language Processing Homework 2

June 6, 2019



**SAPIENZA**  
UNIVERSITÀ DI ROMA

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Model and Approach</b>	<b>2</b>
2.1	Dataset and tools . . . . .	2
2.2	Parsing . . . . .	2
2.3	Training . . . . .	2
2.3.1	Scoring: Similarity measure . . . . .	2
2.3.2	Hyperparameters . . . . .	3
<b>3</b>	<b>Results</b>	<b>3</b>
<b>4</b>	<b>Discussion</b>	<b>4</b>
4.1	General . . . . .	4
4.2	Embedding visualization . . . . .	4
<b>5</b>	<b>Conclusion</b>	<b>5</b>

# 1 Introduction

This report tries to present the process of learning sense embeddings from the annotated EuroSense corpus to be used on a word similarity task.

## 2 Model and Approach

### 2.1 Dataset and tools

The dataset used was the high precision version of the Eurosense corpus [1]. It consists of a large XML file which consists of many multi-lingual sentences taken from the European Union parliament. All the processing steps used in this project were performed in Python 3, with the help of key external libraries. These were **Gensim** used for model training, **iterparse** used for parsing the large corpus, and **Scikit-Learn** for dimensionality. The code file structure can be found in the **README.md** file of the Gitlab repository [4].

### 2.2 Parsing

The parsing consisted of iteratively reading the large corpus file, and processing it line by line. Using special tags, sentences were extracted along with their annotations. However, for the purposes of this assignment extraction was performed only for the English language. Each *sentence*, consisted of several *texts* for the English Language. However, not every *text* had corresponding anchored annotations. Hence, these texts not containing any *anchors* were discarded. Conversely, due to inconsistencies in the dataset, many sentences which had English *anchor* annotations but lacked *texts* were also discarded. Moreover annotations in this corpus file contain a certain ID called *synsetID* which was specified in BabelNet. In this task, only annotations which had a corresponding *synsetID* in WordNet were considered. Stop words were **not** removed, but certain punctuation marks were as they conflicted with some anchors. For example: "European Union-Ukraine Agreement", where there was an anchor for "European Union". Additionally prior to writing the parsed sentences to file, they were lowered as not to have duplicates of senses i.e.: "lemma\_synsetID" and "Lemma\_synsetID". In total, 1.87 million sentences were extracted from the dataset. Figure 1 portrays a histogram of annotation occurrences in each parsed sentence. The median number of annotations was 5. This could be an indicator of how word senses affect each other depending on the number of annotations. In total, 1.87 million sentences were extracted from the dataset.

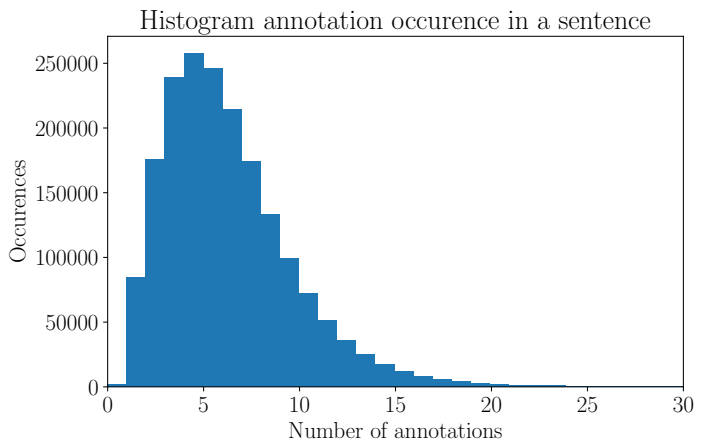


Figure 1: Histogram based on the number of annotations in each of the 1.87 m extracted sentences

### 2.3 Training

The training was performed using the the Word2Vec model from **gensim**. At first, only the *CBOW* version was used, but at the end *SKIPGRAM* was implemented for comparison.

#### 2.3.1 Scoring: Similarity measure

A key target of training the model, was to evaluate it on a similarity measurement task. For this, I utilised cosine similarity between words. Then using Spearman correlation, I tested the results against a gold, human-annotated dataset. I did not incorporate any other similarity measures.

### 2.3.2 Hyperparameters

The key hyperparameters explored when training the model were:  $\alpha$  – learning rate, min count – minimum count of word occurrence in the vocabulary (of the parsed dataset), window, and the embedding size. These hyperparameters were explored using grid search in order to maximize the correlation score. As for the grid search, I implemented a `Word2VecModel` class, which resembles a `search` with best ones in bold `scikit-learn` estimator. It consists of a `fit` function which builds the model, given a set of hyperparameters, then builds the vocabulary, and proceeds to train it on a set number of epochs (50). Then using `sklearn.model.selection.ParameterGrid`, I iteratively perform a grid search of the best hyperparameters based on the cosine similarity score performance.

window	4	<b>5</b>
min count	1	<b>3</b>
embedding size	<b>100</b>	300
$\alpha$	0.001	<b>0.009</b>
Negative	<b>13</b>	
Epochs	<b>50</b>	

## 3 Results

Since *CBOW* works several times faster to train than skipgram, it was used for the grid search. Grid search was first performed on large number of combinations of hyperparameters, on 15 epochs. However, after noticing significant discrepancies in the correlation score, a selected few were selected to be trained in the final grid search. The hyperparameters are presented in Table 1, and the ones in bold refer to the best result of correlation = 0.27 (using *CBOW*). Using these best hyperparameters, a *SKIPGRAM* version of the word2vec model was also trained and resulted in a slightly worse of correlation = 0.26. The difference in performance seems consistent with the fact that *CBOW* has slightly better accuracy for the frequent words, since we used a min count of 3 but *SKIPGRAM* works better with rare phrases. A more graphical approach to visualizing the effect of hyperparameters on each other, and the ones which have the most effect on the final score is presented in Figure 2. The figure presents a swarm plot of the correlation score in regards to the hyperparameters.

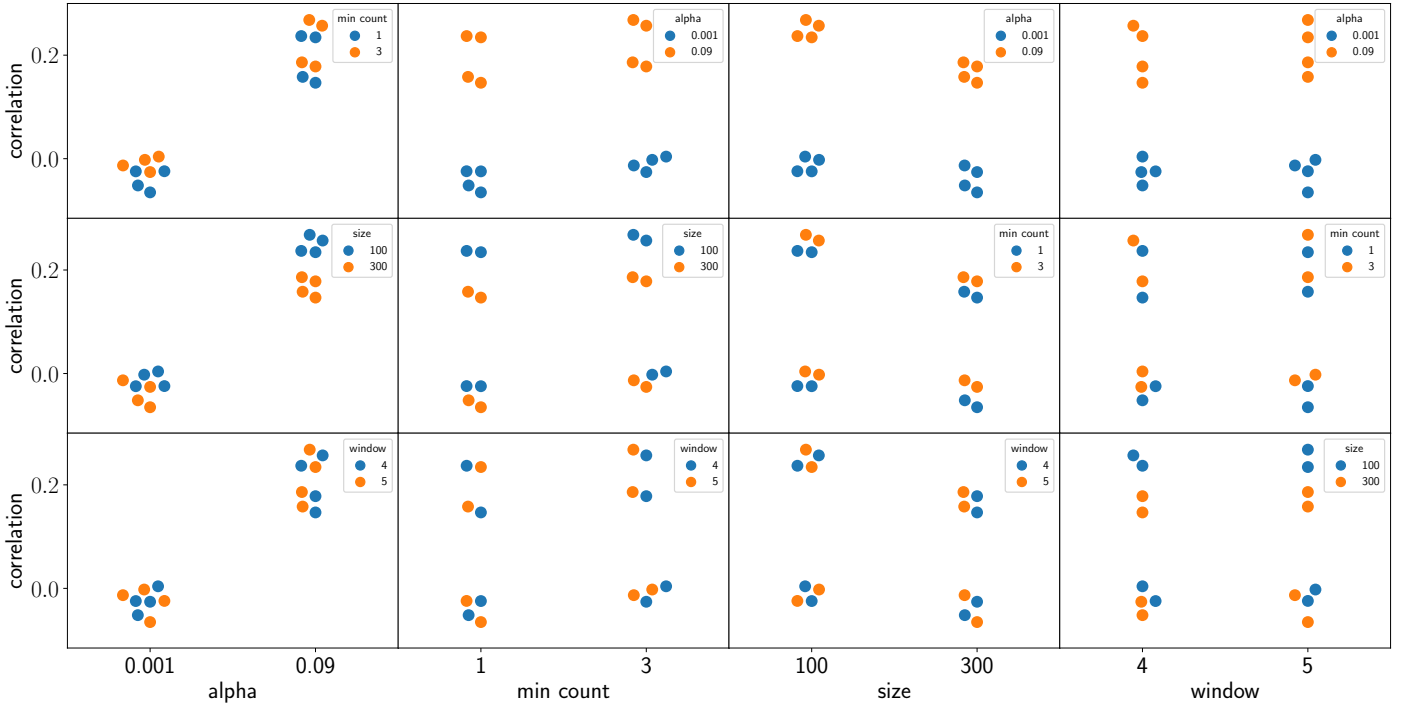


Figure 2: Grid search of the best hyperparameters based on Table 1. Visualization of the Correlation score vs each parameter depending on several others.

## 4 Discussion

### 4.1 General

From Figure 2, we can infer that alpha  $\alpha$ , has the most significant impact on the final correlation score, as regardless of the comparing hyperparameter,  $\alpha = 0.09$  always yields a better correlation score. Moreover, the difference is so substantial that for  $\alpha = 0.001$  correlation scores swarm around 0, which would infer almost the lack of correlation. Additionally, we can see that the size of the embedding also plays a major role, as smaller (100) embeddings perform much better than larger (300). This is a consequence of the size of the vocabulary which is effected both by the min count parameter, but also more importantly by the contents of the corpus.

The corpus is very domain specific which influences correlation scores, with its context and also results in many out-of-vocabulary (OOV) words. A possible solution would be to simply add data from a wider domain, including non factual texts like novels. This however is not easy, as finding enough annotated data is quite difficult but could be partially overcome by combining the already made synsets and analysing them in different ways. On the other hand, it would be far safer to use a good corpus for a more wider domain, i.e: the Semantically Enriched Wikipedia (SEW) [2] corpus.

Furthermore, the overall performance of the model could be yet improved by considering the rising problem of sparsity in the data. One key driving factor of sparsity is the reduction of polysemy which yields a large number of synsets. To overcome this, one could possibly reduce the number of senses for a given lemma but different synset.

### 4.2 Embedding visualization

To further understand what kind of embeddings `word2vec` produces, embedding visualizations were performed. To do this, I reduced their dimensionality from `size = 100` to 2D and 3D using the `sklearn` implementation of Principal Component Analysis (PCA) followed by T-distributed Stochastic Neighbor Embedding (t-SNE) and removed synset IDs for plotting purposes.

Figure 3 depicts a 2D plot of a few chosen word embeddings (legend) and their top 30 most similar words according to the model. From the depicted 4 words we can deduce that they are some are more closely matched than others. Whereas Figure 4 portrays a 3D visualization of all the 44 thousand words in the vocabulary of the model. The plot is basically just one giant cluster in 3D symbolizing that on a larger scale, how much of a domain-specific dataset Eurosense really is.

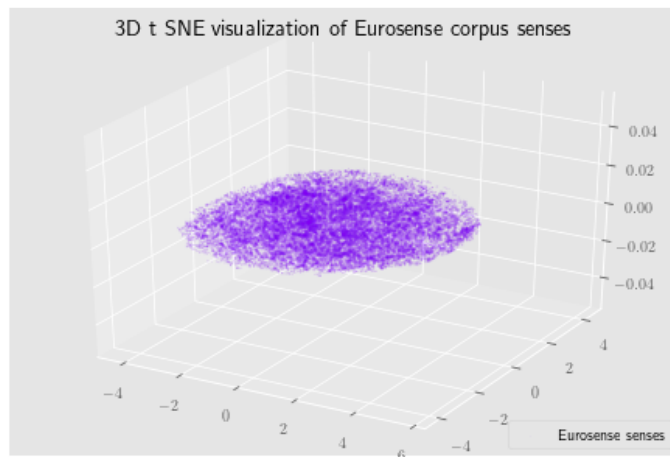


Figure 4: 3D T-SNE on the entire vocabulary of the model. *Perplexity* = 3 and *n\_iter* = 300. First reduced by PCA to 10 dimensions then applied T-SNE on 10 nearest neighbours. Plots inspired by [3]

