

Semantic Embeddings in IR

Text-Based Information Retrieval

Joren Verspeelt, Jochem Geussens, Vincent Tanghe

May 4, 2016

1 Introduction

In the field of Natural Language Processing (NLP) and Information Retrieval (IR), there is a large need for good presentations of text documents. Recently, a lot of researchers in the field are presenting dense word representations (also known as Word Embeddings, Neural Embeddings or Semantic Embeddings). For the Course of Text-Based Information Retrieval, we are given the opportunity to work with state-of-the-art algorithms. The goal of this assignment is to gain practical experience with these algorithms by comparing them and using them in an application. First we will implement a small analogy solver that guesses related words given a simple analogy. Then, we will discuss our implementation for a search engine to find pictures based on their description.

Contents

2 Part I (Warm-up)

2.1 Task description

In this part, we discuss two analogy solvers that try to guess the correct word, given simple analogy questions. E.g.:

$$\begin{aligned} a : b &= c : ? \\ Bratislava : Slovakia &= Bishkek : ? \\ ate : eat &= found : ? \end{aligned}$$

Then, we will run several analogy solving models with several different representations on the benchmarking analogy dataset.

2.2 Setup

We compare the use of different word embeddings (GloVe and Word2Vec) towards different analogy arithmetic models[1]. To test these, we used the questions provided by Google[2].

We coded our solution in Python (version 2.7) using IPython Notebook. All experiments were run on a notebook with 8GB RAM memory, Intel(R) CORE™i7-3610QM CPU, SSD hard disk and a Windows 7 Operating System.

We use freely available pretrained vector models for GloVe [3] and Word2Vec [4]. The GloVe models are almost the same format as the Word2Vec models. The only difference is that Word2Vec had a header with the dimensions. Once added, the GloVe model can also run with the default Gensim operations for Word2Vec.

To run the different analogy arithmetic models (see next paragraph), we introduce a separate function in our code for each model. These functions then use the Gensim library for Python [5] to calculate the different vector distances.

An analogy arithmetic function is used to calculate how close the input vector is to another vector in order to choose the most similar result. To calculate the angle between high dimensional vectors, Mikolov et al. [6] proposed using cosine similarity as explained in equation 1.

$$\begin{aligned} similarity(v1, v2) &= \cos(\theta) \\ &= \frac{v1 \cdot v2}{||v1|| \cdot ||v2||} \\ &= \frac{\sum_{i=1}^n v1_i v2_i}{\sqrt{\sum_{i=1}^n v1_i^2} \sqrt{\sum_{i=1}^n v2_i^2}} \end{aligned} \tag{1}$$

We consider two analogy arithmetic models. The first is named the **Addition method** and is simply maximizing the cosine similarity between the result d_w and the vector equation $c - a + b$, where the variables a, b and c have the same meaning as mentioned in the warm-up exercise. This method actually looks for words similar to c and b , and not to a . In equation ?? we introduce V as the whole vocabulary and d'_w as a chosen element from V .

$$\begin{aligned} v &= c - a + b \\ d_w &= \operatorname{argmax}_{d'_w \in V} (\operatorname{similarity}(d'_w, v)) \end{aligned} \tag{2}$$

The other method was referred to as the **Multiplication method** [7]. Because the Addition method risks being dominated by one large term, therefore the authors proposed that instead of adding the similarities, they could be multiplied as described in equation 3.

$$d_w = \operatorname{argmax}_{d'_w \in V} \left(\frac{\operatorname{similarity}(d'_w, c) \operatorname{similarity}(d'_w, b)}{\operatorname{similarity}(d'_w, a) + \epsilon} \right) \tag{3}$$

Two different recall numbers are generated if a word does not occur within the pretrained vector model: One where the **missing word is considered a failure** for the analogy and one where the **missing word is just ignored**. We could also have done a third option, were we use the nearest word in the vector model instead of the missing word. However, we did not do this because this would take a lot more computing power and (as seen in the results below) we're not sure whether it would have made a difference.

2.3 Results

In table 1 the results of the GloVe method combined with the multiplication method are shown. The results for the addition method are displayed in table ??. The results of multiplication- and addition model of Word2Vec can respectively be found in table 3a and table 3b. A graph that visualizes the results next to each other is presented in figure 1.

No matter the dimension, GloVe always had the same amount of skipped words (all words) within a category. So we will not show the category results for these.

Category	Recall
Family	0.62846
Gram1 adjective to adverb	0.09980
Gram2 opposite	0.04926
Gram3 comparative	0.41366
Gram4 superlative	0.17112
Gram5 present participle	0.29451
Gram7 past tense	0.31153
Gram8 plural	0.46021
Gram9 plural verbs	0.25862
Total	0.14506
Total no skipped	0.29587

(a) GloVe50d multiplication model (3 min run-time)

Category	Recall
Family	0.77866
Gram1 adjective to adverb	0.22883
Gram2 opposite	0.15764
Gram3 comparative	0.74850
Gram4 superlative	0.50624
Gram5 present participle	0.65057
Gram7 past tense	0.52821
Gram8 plural	0.66667
Gram9 plural verbs	0.57356
Total	0.26668
Total no skipped	0.54394

(b) GloVe100d multiplication model (7 min run-time)

Category	Recall
Family	0.85178
Gram1 adjective to adverb	0.25302
Gram2 opposite	0.19581
Gram3 comparative	0.83784
Gram4 superlative	0.67162
Gram5 present participle	0.67045
Gram7 past tense	0.60321
Gram8 plural	0.76426
Gram9 plural verbs	0.65172
Total	0.03531
Total no skipped	0.62273

(c) GloVe200d multiplication model (10 min run-time)

Category	Recall
Total	0.32214
Total no skipped	0.65707

(d) GloVe300d multiplication model (20 min runtime)

Table 1: Results of GloVe multiplication models

Category	Recall
Family	0.85573
Gram1 adjective to adverb	0.25403
Gram2 opposite	0.22660
Gram3 comparative	0.86486
Gram4 superlative	0.69786
Gram5 present participle	0.68277
Gram7 past tense	0.60192
Gram8 plural	0.77402
Gram9 plural verbs	0.59770
Total	0.30778
Total no skipped	0.62774

(a) GloVe50d addition model (3 min runtime)

Category	Recall
Family	0.81621
Gram1 adjective to adverb	0.24395
Gram2 opposite	0.20074
Gram3 comparative	0.79129
Gram4 superlative	0.54278
Gram5 present participle	0.69508
Gram7 past tense	0.55448
Gram8 plural	0.71997
Gram9 plural verbs	0.58391
Total	0.28382
Total no skipped	0.57890

(b) GloVe100d addition model (7 min runtime)

Category	Recall
Family	0.85573
Gram1 adjective to adverb	0.25403
Gram2 opposite	0.22660
Gram3 comparative	0.86486
Gram4 superlative	0.69786
Gram5 present participle	0.68277
Gram7 past tense	0.60192
Gram8 plural	0.77402
Gram9 plural verbs	0.59770
Total	0.30778
Total no skipped	0.62774

(c) GloVe200d addition model (10 min runtime)

Category	Recall
Total	0.31304
Total no skipped	0.63849

(d) GloVe300d addition model (20 min runtime)

Table 2: Results of GloVe addition models

Category	Recall
Capital common countries	0.85178
Capital World	0.80570
Currency	0.35450
City in state	0.71423
Family	0.84585
Gram1 adjective to adverb	0.31552
Gram2 opposite	0.42365
Gram3 comparative	0.90691
Gram4 superlative	0.91800
Gram5 present participle	0.80587
Gram6 nationality adjective	0.89368
Gram7 past tense	0.70641
Gram8 plural	0.89940
Gram9 plural verbs	0.73448
Total	0.75148

(a) Word2Vec multiplication model (3h25m runtime)

Category	Recall
Capital common countries	0.83202
Capital World	0.79134
Currency	0.35104
City in state	0.70896
Family	0.84585
Gram1 adjective to adverb	0.28528
Gram2 opposite	0.42734
Gram3 comparative	0.90841
Gram4 superlative	0.87344
Gram5 present participle	0.78125
Gram6 nationality adjective	0.89931
Gram7 past tense	0.65962
Gram8 plural	0.89865
Gram9 plural verbs	0.67931
Total	0.73588

(b) Word2Vec addition model (1h15m runtime)

Table 3: Results of Word2Vec models

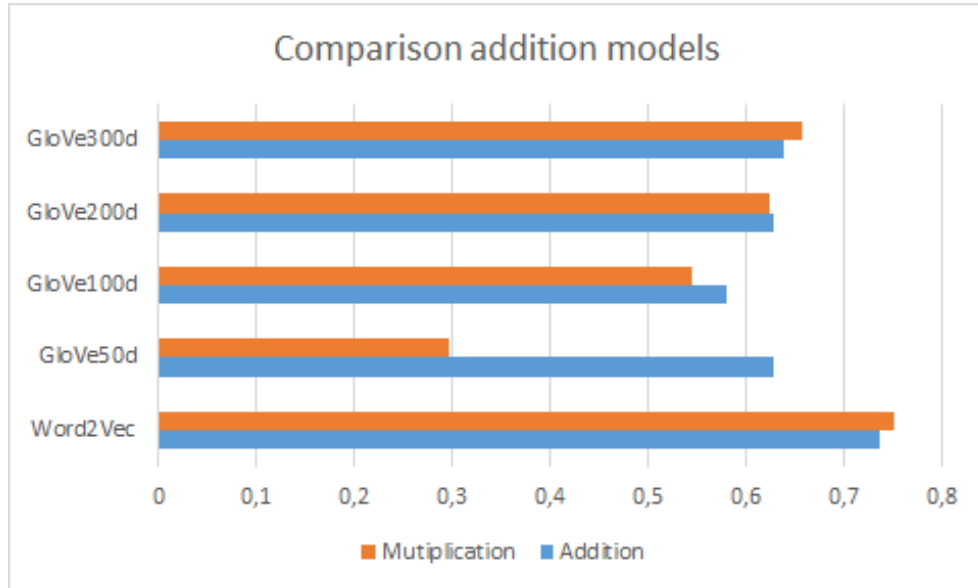


Fig. 1: A comparison between GloVe and Word2Vec results for the multiplication and addition method.

2.4 Discussion

We noticed that the GloVe model was not able to find all the words and this influences the results regarding execution times. The GloVe model seems to run a lot faster than the Word2Vec model because of our implementation strategy. If one word of the analogy is missing, the whole analogy is skipped resulting in smaller search space that needs to be covered. But that aside the GloVe model has the edge with faster calculations due to its lower dimensionality.

If we look at the dimensionality of the GloVe representations, we see that the overall accuracy increases as the number of dimensions increases. This is a trade-off compared to the time needed to run the model, but in general it seems worthwhile. If we then look at the analogy model compared to the dimensionality, we see that the multiplication model performs better on higher dimensions while the addition model performs better as the dimension decrease. The overall accuracy decrease with a lower dimension is most likely the reason for this.

While running the experiment, we frequently got missing words (using the GloVe representation). It is still impractical to have a model representing every possible word that exists, so a trade-off is made towards a reasonable model size and the amount of represented words. We solved this problem by either counting missing words as failed analogies or by skipping the sentence. However, it seems that the same word categories are missing in every dimension.

3 Part II (Retrieval task)

3.1 Task description

This task is based on the annual Scalable Concept Image Annotation Challenge, organised by ImageCLEF [8]. The dataset (given by ImageCLEF), contains sentence descriptions aligned to actual images. In this part of the assignment, we bring our previous tests into a real-life retrieval scenario in order to retrieve image descriptions using textual queries. Each query has a corresponding image ID, which can be used to verify the correctness of the model. The goal is to match queries with image descriptions as good as possible, using precision and recall as evaluation measures.

3.2 Setup

We coded our solution in Python (version 2.7) using IPython Notebook. All experiments were run on a notebook with 8GB RAM memory, Intel(R) CORE™i7-3610QM CPU, SSD hard disk and a Windows 7 Operating System.

The data we received contains the following documents:

- A. target_collection** A table representing 17,784 images with as columns a unique numerical ID, an ID that corresponds to the actual image file and a description of the image. This data is used to train our model.
- B. queries_val** A similar table as the *target_collection* but it represents 1000 different images. This table will be used to test the accuracy of our trained model.

3.3 Results

3.3.1 Unigram Language Model

3.3.2 title

3.4 Discussion

4 Conclusion

[Strong points, Weak points, lessons learned]

5 References

- [1] M. Rei, “*Linguistic regularities in word representations*,” 2014. <http://www.marekrei.com/blog/linguistic-regularities-word-representations> [Accessed: 3-5-2016].
- [2] “*Question words provided by Google Code*,” 2013. <http://word2vec.googlecode.com/svn/trunk/questions-words.txt> [Accessed: 3-5-2016].
- [3] C. D. M. Jeffrey Pennington, Richard Socher, “*GloVe: Global Vectors for Word Representation*,” 2014. <http://nlp.stanford.edu/projects/glove/> [Accessed: 3-5-2016].
- [4] “*Word2vec*,” 2013. <https://code.google.com/archive/p/word2vec/> [Accessed: 3-5-2016].
- [5] R. Rehurek, “*gensim: topic modelling for humans*,” 2016. <https://radimrehurek.com/gensim/> [Accessed: 3-5-2016].
- [6] Y. W. Z. G. Mikolov, T., *Linguistic Regularities in Continuous Space Word Representations*, 746–751.
- [7] G. Y. Levy, O., *Linguistic Regularities in Sparse and Explicit Word Representations. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (pp. 171–180)*.
- [8] “<http://imageclef.org/2016/annotation>.” <http://imageclef.org/2016/annotation>.