

# WORD EMBEDDINGS FROM LARGE-SCALE GREEK WEB CONTENT

*Stamatis Outsios<sup>1</sup>, Konstantinos Skianis<sup>2</sup>, Polykarpos Meladianos<sup>1</sup>,  
Christos Xypolopoulos<sup>2</sup>, Michalis Vazirgiannis<sup>1,2</sup>*

<sup>1</sup> Athens University of Economics and Business  
Department of Informatics  
Greece

<sup>2</sup> École Polytechnique  
Laboratoire d'informatique (LIX)  
France

## ABSTRACT

Word embeddings are undoubtedly very useful components in many NLP tasks. In this paper, we present word embeddings and other linguistic resources trained on the largest to date digital Greek language corpus. We also present a live web tool for testing the Greek word embeddings<sup>1</sup>, by offering “analogy”, “similarity score” and “most similar words” functions. Through our explorer, one could interact with the Greek word vectors.

**Index Terms**— Greek word embeddings, Greek web, Greek language NLP resources

## 1. INTRODUCTION & RELATED WORK

With the rise of neural networks and deep learning in the NLP community [1, 2], word embeddings were introduced [3], having a huge impact on numerous tasks. Their ability to represent rich relationships between words, led to state-of-the-art results, combined with CNNs [4] and LSTMs [5] in text classification, question answering and machine translation tasks. Although word vector transformation is the most common way to harvest text, they require a large amount of data in order to be trained. Moreover, resources for specific languages may be scarce or hard to extract in an efficient way.

In this work, we present Greek word embeddings, trained on, to the best of our knowledge, the largest so far corpus available, collected/crawled from about 20M URLs with Greek language content. The vocabulary and word vectors are available on request. We developed a live web tool to enable users to interact with the Greek word embeddings. Some of the functions we provide is similarity score, most similar words as well as analogy operations. We also present a vector explorer, where we can project a sample of the word vectors.

Lately, pre-trained word vectors for 294 languages were introduced, trained on Wikipedia using FastText[6]. These 300-dimensional vectors were obtained using the skip-gram model. Their Greek language variant, was trained on the Wikipedia corpus concerning only Greek documents.

Visualization tools for word embeddings are of great importance, since they contribute to the interpretation of their nature. Similarly to our tool, Tensorflow<sup>2</sup> offers an illustration of a sample of word embeddings after applying dimensionality reduction techniques. Last, the training process can be observed with WEVI, a word embedding visual inspector<sup>3</sup>.

## 2. CRAWLING THE GREEK WEB

For the process of crawling the Greek Web (which was funded by the Stavros Niarchos foundation, see: <https://www.snf.org>, for the Greek National Library) we used the Heritrix<sup>4</sup> tool. Collecting the websites adheres to the international Web Archive (WARC) template. The WARC file form defines a method combining multiple media resources into one archive. Next, we present some statistics about the data we crawled:

- Number of WARCs: 112K
- Size of HTML (stored in WARC format): 10TB
- Number of Greek domains: 350K
- Number of URLs: 20M
- Duration of crawling: 45 days

## 3. PRE-PROCESSING & TEXT EXTRACTION

Before training, we applied several pre-processing and extraction steps on the raw crawled text:

1. detect the encoding of each webpage, so that we are able to read it properly,
2. remove HTML code and tags, as well as Javascript,
3. remove boilerplate code<sup>5</sup>,
4. remove all non-Greek characters,
5. track the line change character,
6. produce compressed text files per domain.

The third step is very important for the text quality, since we request a corpus that can be used later for developing linguistic resources (language model, embeddings etc.). Except their

<sup>1</sup><http://archive.aueb.gr:7000>

<sup>2</sup><https://projector.tensorflow.org/>

<sup>3</sup><https://ronxin.github.io/wevi/>

<sup>4</sup><http://crawler.archive.org/>

<sup>5</sup>[https://en.wikipedia.org/wiki/Boilerplate\\_code](https://en.wikipedia.org/wiki/Boilerplate_code)

content, webpages consist of navigation elements, headers, footers, as well as commercial banners. This text is usually not associated with the webpage’s main content, and can lead in decreasing the integrity of the collection. In order to do that, we used libraries like BeautifulSoup, Justext, NTLKs clean\_html and Boilerpipe<sup>6</sup>. The best results were obtained by Boilerpipe, which was the one we used in the end to remove useless text (boilerplate). We removed identical sentences (de-duplication) and produced the final corpus in text form, sized around 50GB. We obtained thus ~3B tokens and a total number of 498M sentences, with 118M of them being unique.

De-duplication per domain resulted in reducing the size of raw corpus by 75%. With an additional processing of the final corpus, we create the Greek language n-grams: Unigrams: ~7M, Bigrams: ~90M, Trigrams: ~300M.

#### 4. TRAINING GREEK WORD EMBEDDINGS

For the process of learning the Greek word embeddings we utilized the FastText[6] library, which takes under consideration the morphology of a word. Training on the raw uncompressed text of the Greek internet web, with size of 50GB, required 2 days in a 8-core Ubuntu system with 32GB of RAM.

Different Greek vector models were produced like: 1. native fasttext skipgram with the following parameters: -minCount 11 -loss hs -thread 8 -dim 300, 2. native fasttext cbow, 3. gensim<sup>7</sup> word2vec skipgram, 4. gensim fasttext skipgram, 5. gensim fasttext skipgram, no subword information. Methods 3 and 4 lead to the same result as they use the same technique. By evaluating their effectiveness in automatic correction along with similarity queries, method 1 yields the most reliable results. In the future, we plan to offer as well a set of handcrafted questions for evaluation purposes.

#### 5. VISUALIZATION

Next, we designed tools in order to visualize examples of Greek word vector relationships. The first demo offers linguistic functions which are enabled by the existence of word embeddings, like analogy, similarity score or most similar words. The second demo tool for exploring and querying the word vectors was based on the word2vec-explorer<sup>8</sup>. In this tool, a user can navigate through a sample of the Greek word embeddings, visualize it via t-SNE[7] and apply k-means clustering. Furthermore, we offer comparing functions for combinations of words. For the frontend, we used libraries like Flask, Jinja and Bootstrap.

<sup>6</sup><https://boilerpipe-web.appspot.com/>

<sup>7</sup><https://radimrehurek.com/gensim/>

<sup>8</sup><https://github.com/dominiek/word2vec-explorer>

#### 6. CONCLUSION & FUTURE WORK

In this work, we present the efforts that resulted in Greek word embeddings and other Greek language resources trained on the largest corpus available the Greek language. The resources (corpus, trained vectors, stopwords, vocabulary as well as unigrams, bigrams and trigrams) are available on request. We have also implemented a live web tool, where a user can explore word relationships in the Greek language. In addition, we provide a word embedding explorer, where one could visualize a sample of the Greek vectors with t-SNE[7].

Recently, embeddings evolved towards new approaches like Hierarchical Representations [8] or ELMo [9]. Finally, we plan to extend our work by adding a visualization of the Word Mover’s Distance[10].

#### 7. REFERENCES

- [1] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin, “A neural probabilistic language model,” *JMLR*, 2003.
- [2] Ronan Collobert and Jason Weston, “A unified architecture for natural language processing: Deep neural networks with multitask learning,” in *ICML*, 2008.
- [3] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean, “Distributed representations of words and phrases and their compositionality,” in *NIPS*, 2013, pp. 3111–3119.
- [4] Yoon Kim, “Convolutional neural networks for sentence classification,” in *EMNLP*, 2014.
- [5] Rie Johnson and Tong Zhang, “Supervised and semi-supervised text categorization using lstm for region embeddings,” in *ICML*, 2016.
- [6] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov, “Enriching word vectors with subword information,” *TACL*, 2017.
- [7] Laurens van der Maaten and Geoffrey Hinton, “Visualizing data using t-sne,” *JMLR*, 2008.
- [8] Maximilian Nickel and Douwe Kiela, “Poincaré embeddings for learning hierarchical representations,” in *NIPS*. 2017.
- [9] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer, “Deep contextualized word representations,” in *NAACL*, 2018.
- [10] Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger, “From word embeddings to document distances,” in *ICML*, 2015.