

# Optimized Training and Evaluation of Arabic Word Embeddings

Jordan King, Lisa Singh, Eric Wang

November 1, 2015

## Abstract

Word embeddings are an increasingly important tool for NLP tasks, especially those that require semantic understanding of words. Methodologies and properties of word embeddings have been researched in English language tasks, and tasks with English-like languages. However, little attention has been given to the training of Arabic word embeddings. Arabic is much more morphologically complex than English, and due to the many conjugations, suffixes, articles, and more, an Arabic word might not be a single word in the sense that it is in English. While there are a number of techniques to break down Arabic words through lemmatization and tokenization, it is not clear how the quality of resulting word embeddings would be affected. We investigate a number of preprocessing methods and training parameterizations to find the optimal strategies to train embeddings with. Additionally, we required tasks to evaluate the Arabic word embeddings. There is little work done in providing a semantic similarity task for Arabic. To remedy this, we created a list of human supplied similarity scores for given Arabic word pairs. We evaluated all training methods by using the embeddings from a given method to obtain a similarity score, and compare that against our semantic similarity task. Additionally, we evaluated the word embeddings ability to capture syntactic properties using a part of speech tagging task. Using these tasks, we are able to identify the training strategies that perform best on each task. We also offer a suite of Arabic NLP tools that we developed alongside this work that attempts to fill the void of accessible open source Arabic NLP tools. Altogether, this work provides best practices for training Arabic word vectors, an open semantic simi-

larity task developed by native Arabic speakers, and a python package of Arabic NLP tools.

## 1 Introduction

We are researching methods to create Arabic word embeddings, which are numerical representations of a given words meaning - both semantic meaning and syntactic meaning. These embeddings are obtained using machine learning algorithms that utilize the context a word appears in to infer its meaning. This works very well as words with similar meanings tend to be used in similar contexts. For example "I eat bread every night" and "I eat rice every night" is an example of how foods may appear in similar contexts. With enough text to process, we can train numerical vectors to learn that bread and rice appear in these "common-for-food" contexts. Similarly, we can learn syntactic relationships as different parts of speech appear in certain context patterns as well.

The end result of having good word embeddings is that we have a representation of the meaning of a word, without ever translating or looking at a dictionary of meanings. We can harvest the semantic and syntactic meaning straight from a corpus of natural written language. With accurate word embeddings, we can perform neat and useful operations to investigate the relationships between words. To list a few of these operations, we can measure the similarity of two words, identify which word from a set is least similar, and solve basic analogies. The classic party trick for word embeddings is to take (the numerically embeddings of) king, subtract man, and add woman. The resulting embedding is closest to the embedding for queen! Intuitively, this allows us to subtract the male gender meaning from kings embedding, add the female gender meaning, and end up with an embedding equivalent to queens embedding.

More about the algorithms and uses: <https://code.google.com/p/word2vec/>

We would like accurate Arabic word embeddings so we can interpret the general topics of discussion in Arabic media without using translation or ignoring some words belonging to a topic. Specifically, one task we are using the embeddings for is to learn what words are highly similar to words such as fear (in Arabic), and then compute the degree to which some media is using these words similar to fear in the context of some city. This information may help us predict when people will be forced to migrate from the city.

**2 Related Literature**

**3 Experiments**

**4 Conclusion**