



INM713 Semantic Web Technologies and Knowledge Graphs

Laboratory 9: Ontology Embeddings

Ernesto Jiménez-Ruiz

Academic course: 2020-2021

Updated: March 22, 2021

Contents

1	Git Repositories	2
2	Using OWL2Vec*	2
3	Using the Embeddings	4
4	Embedding the Pizza ontology	4
5	Embedding the FoodOn ontology	5

1 Git Repositories

Support codes for the laboratory sessions are available in *GitHub*. There are two repositories, one in Python and another in Java:

```
https://github.com/city-knowledge-graphs
```

2 Using OWL2Vec*

In this laboratory session we are using OWL2Vec*, a system that embeds the semantics of an OWL ontology (<https://github.com/KRR-Oxford/OWL2Vec-Star>). OWL2Vec* currently relies on random walks and word embedding and learns a (numerical) vector representation for each URI and word in the ontology. OWL2Vec* generates a document of sentences from the ontology (using this RDF2vec implementation: <https://github.com/IBCNServices/pyRDF2Vec/>) and then applies Word2vec as neural language model (see details here: <https://radimrehurek.com/gensim/models/word2vec.html>).

The codes for OWL2Vec* are available in the GitHub repositories as a zip file `owl2vec_standalone.zip`. Unfortunately for the Java developers, OWL2Vec is only available in Python. Nevertheless, running OWL2Vec*, once the dependencies have been installed (see below), is as easy as running the following command:

```
python OWL2Vec_Standalone.py -config_file default.cfg
```

In `default.cfg` we can indicate the following parameters (see lecture slides or OWL2Vec* paper for more details):

- `ontology_file`: input OWL ontology.
- `embedding_dir`: path and file name for the output embeddings.
- `cache_dir`: path to store intermediate files.
- `ontology_projection`: if the graph projection of the ontology is enabled.
- **Walker parameters**: to build the document sentences.
 - `walker`: random walks or Weisfeiler-Lehman (wl) subtree kernel.
 - `walk_depth`: depth of each of the performed walks to create each sentence.
- **OWL2Vec* paramters**: type of document of sentences to build. One could create a document with only words (*i.e.*, `Lit_Doc`)
 - `URI_Doc`: sentences with only entity URIs.
 - `Lit_Doc`: sentences with only words.
 - `Mix_Doc`: mixing words and URIs in the sentences.

- `Mix_Type`: the type for generating the mixture document (all or random).
- `pre_train_model`: optional path to a pre-trained Word2vec model.
- **Word2vec parameters:**
 - `embed_size`: the size for embedding.
 - `iteration`: number of iterations in training the language model.
 - `min_count`: minimum word occurrence to create an embedding.
 - `window`, `negative` and `seed`: other Word2vec training parameters.

OWL2Vec* produces three embedding files as output (in `embedding_dir`):

- Gensim model (`.embedding` file).
- Word2vec binary format (`.bin` file).
- Word2vec text format (`.txt` file). This file is readable with a text editor. Each line is composed by a key (*i.e.*, word) and a value (*i.e.*, vector).

The generated sentence document is available in the `cache` output folder (*i.e.*, `document_sentences.txt`). This file is interesting to see how the ontology has been transformed into a text document. It is also very relevant for the Java developers (see below).

Dependencies. OWL2Vec* requires the following additional libraries.

- <https://pypi.org/project/gensim/>
- <https://pypi.org/project/scikit-learn/>
- <https://pypi.org/project/nltk/>
- <https://pypi.org/project/numpy/>

In addition, Java developers will also need to install:¹

- <https://pypi.org/project/Owlready2/>
- <https://pypi.org/project/rdfliib/>

¹Assuming that they have already installed Python 3. Otherwise install from <https://www.python.org/downloads/>

3 Using the Embeddings

The computed embeddings can be used to calculate (cosine) similarities, perform clustering of entities, and use them in subsequent machine learning tasks.

Python. In the GitHub repository there is an example script (`load_embeddings.py`) that loads pre-computed embeddings, in this case by OWL2Vec*. In python we use the `gensim` `keyedvectors` library: <https://radimrehurek.com/gensim/models/keyedvectors.html>. Once the model has been loaded one can calculate (cosine) similarities among ontology entities/words (*i.e.*, using their associated vectors) and get their closest entities/words.

Java. In the GitHub repository, there is a class `WordEmbeddings.java`, implemented over the `Deeplearning4j` library (<https://deeplearning4j.konduit.ai/language-processing/word2vec>). Although it is in principle possible to load in Java embedding models pre-computed in Python, this did not seem to work. `WordEmbeddings.java` creates embeddings from the document of sentences (*i.e.*, `document_sentences.txt`). As in Python, with this class one can also load a pre-computed model and calculate similarities among ontology entities/words (*i.e.*, using their associated vectors) and get their closest entities/words. The results with the java version of `Word2vec` differ with respect to those in Python. Alternatively, one could also try to load the vectors directly reading the `Word2vec` text file computed by OWL2Vec*.

4 Embedding the Pizza ontology

The ontology is available in the GitHub repositories.

Task 4.1 Run OWL2Vec* over the `pizza.owl` ontology.

Task 4.2 Compute the similarity between the following elements.

- <http://www.co-ode.org/ontologies/pizza/pizza.owl#Margherita> and `margherita`
- <http://www.co-ode.org/ontologies/pizza/pizza.owl#Hot> and <http://www.co-ode.org/ontologies/pizza/pizza.owl#Medium>
- `CheesyPizza` and `CheeseTopping`

Task 4.3 Get the most similar entities/words for the following elements:

- `Hot`
- <http://www.co-ode.org/ontologies/pizza/pizza.owl#CheesyPizza>
- <http://www.co-ode.org/ontologies/pizza/pizza.owl#Fiorentina>
- `Soho`

Task 4.4 Perform clustering using the K-means algorithm for example and visually represent the results in 2-dimensions using PCA.

Python, relevant references:

- <https://jakevdp.github.io/PythonDataScienceHandbook/05.11-k-means.html>
- <https://towardsdatascience.com/pca-using-python-scikit-learn-e653f8989e60>

Java, relevant references:

- <https://deeplearning4j.org/api/latest/org/deeplearning4j/clustering/kmeans/KMeansClustering.html>
- <https://deeplearning4j.org/api/latest/org/nd4j/linalg/dimensionalityreduction/PCA.html>

5 Embedding the FoodOn ontology

In the following tasks we will use the FoodOn ontology (<https://foodon.org/>). The ontology is also available in the GitHub repositories.

Task 5.1 Run OWL2Vec* over the foodon-merged.owl ontology. This ontology is larger and computing the embeddings will take much longer. The embeddings will also be richer as the generated document is rather large (>1Gb).

Task 5.2 Compute the similarity between the following elements.

- http://purl.obolibrary.org/obo/FOODON_00002434 (mushroom food product) and mushroom
- http://purl.obolibrary.org/obo/FOODON_00002434 (mushroom food product) and http://purl.obolibrary.org/obo/FOODON_00001287 (mushroom food source)
- http://purl.obolibrary.org/obo/FOODON_03304544 (frozen chicken) and http://purl.obolibrary.org/obo/FOODON_03311876 (baked chicken)

Task 5.3 Get the most similar entities/words for the following elements:

- mushrooms
- chicken
- http://purl.obolibrary.org/obo/FOODON_03411323 (rabbit)
- http://purl.obolibrary.org/obo/FOODON_03411129 (trout and salmon family)

Task 5.4 Perform clustering using the K-means algorithm for example and visually represent the results in 2-dimensions.