



## PhD course on Knowledge Graphs in the era of Large Language Models

# Ontology Embeddings with OWL2Vec\*

Ernesto Jiménez-Ruiz

April 2024

Updated: April 26, 2024

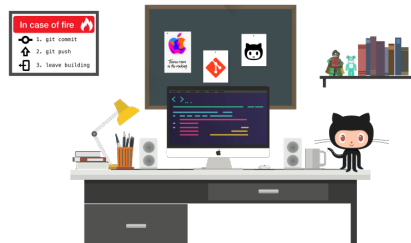
# Contents

<b>1</b>	<b>Git Repository</b>	<b>2</b>
<b>2</b>	<b>Using OWL2Vec*</b>	<b>2</b>
2.1	Installation . . . . .	2
2.2	Configuration . . . . .	3
2.3	Output . . . . .	4
<b>3</b>	<b>Using the Embeddings</b>	<b>4</b>
<b>4</b>	<b>Embedding the Pizza ontology</b>	<b>4</b>
<b>5</b>	<b>Embedding the FoodOn ontology</b>	<b>5</b>

# 1 Git Repository

Support codes for the laboratory sessions are available in *github*.

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



## 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>). OWL2Vec\* has been implemented in Python.

### 2.1 Installation

Please use the the OWL2Vec\* version (zip file) from the GitHub repositories (OWL2Vec-Star-Teaching-2024.zip). Running OWL2Vec\* is very easy. Before you start, download OWL2Vec\* (OWL2Vec-Star-Teaching-2024.zip), unzip the file, and open the folder where the codes are (e.g., location of `setup.py`).

**Tip:** To run the commands in Options 1 and 2 below make sure you are located in that directory, also to run the jupyter notebook that is provided within the zip file.

**Option 1** (running OWL2Vec\* from the terminal):

1. Install Python 3 (if not done before): <https://www.python.org/downloads/>
2. Install setuptools: <https://pypi.org/project/setuptools/>
3. Run this command in the terminal:<sup>1</sup> `make install` or `python setup.py install` (in Windows and other distributions). You may need Root privileges in Linux.

---

<sup>1</sup>For Windows user you need to use the Windows Command Prompt. I can help on this: <https://www.makeuseof.com/tag/a-beginners-guide-to-the-windows-command-line/>

4. Run OWL2Vec\* in the terminal:

```
owl2vec_star standalone --config_file default.cfg
```

**Option 2** (running OWL2Vec\* from a notebook or a python script):

1. Install Python 3 (if not done before): <https://www.python.org/downloads/>
2. Install pip (if not done before): (<https://pip.pypa.io/en/stable/installation/>)
3. Install the library dependencies in the file `requirements_owl2vec.txt` (use of environments is recommended):  
`e.g., pip3 install -r requirements_owl2vec.txt`
4. Run the notebook: `jupyter_notebook_owl2vec.ipynb`

## 2.2 Configuration

In the file `default.cfg` we can indicate the following configuration parameters for OWL2Vec\* (see lecture slides or the 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.
- `projection_only_taxonomy`: if only `rdfs:subClassOf` is taken into account.
- `multiple_labels`: if more than the main label is considered.
- `avoid_owl_constructs`: if OWL 2 constructs are avoided in the generated document.
- **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\* parameters**: 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.

## 2.3 Output

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

- Gensim model (`.embedding` 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.

## 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. A jupyter notebook is also provided.

## 4 Embedding the Pizza ontology

The ontology is available in the GitHub repositories.

**Task 1** Run OWL2Vec\* over the `pizza.owl` ontology.

**Task 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 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** (Optional) 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://jakevdp.github.io/PythonDataScienceHandbook/05.09-principal-component-analysis.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** 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 6** Compute the similarity between the following elements.

- [http://purl.obolibrary.org/obo/FOODON\\_00002434](http://purl.obolibrary.org/obo/FOODON_00002434) (mushroom food product) and mushroom
- [http://purl.obolibrary.org/obo/FOODON\\_00002434](http://purl.obolibrary.org/obo/FOODON_00002434) (mushroom food product) and [http://purl.obolibrary.org/obo/FOODON\\_00001287](http://purl.obolibrary.org/obo/FOODON_00001287) (mushroom food source)
- [http://purl.obolibrary.org/obo/FOODON\\_03304544](http://purl.obolibrary.org/obo/FOODON_03304544) (frozen chicken) and [http://purl.obolibrary.org/obo/FOODON\\_03311876](http://purl.obolibrary.org/obo/FOODON_03311876) (baked chicken)

**Task 7** Get the most similar entities/words for the following elements:

- mushrooms

- chicken
- [http://purl.obolibrary.org/obo/FOODON\\_03411323](http://purl.obolibrary.org/obo/FOODON_03411323) (rabbit)
- [http://purl.obolibrary.org/obo/FOODON\\_03411129](http://purl.obolibrary.org/obo/FOODON_03411129) (trout and salmon family)

**Task 8** (Optional) Perform clustering using the K-means algorithm for example and visually represent the results in 2-dimensions.