

Announcement

21.06.2019

---

MASTER THESIS: **Learning and Evaluating the Quality of Language models.**

## DESCRIPTION

Language Models provides context to distinguish between words and phrases that sound similar used in speech recognition, machine translation, part-of-speech tagging, parsing, Optical Character Recognition, handwriting recognition, information retrieval and other applications.

One example of **Language Model** is Word Embeddings, that is a groundwork for most tasks in the field of Natural Language Processing (NLP). Using Word Embedding representation is possible to capture multiple different degrees of similarity between words. Mikolov et al. (2013) found that semantic and syntactic patterns can be reproduced using vector arithmetic. For instance, patterns such as “Man is to Woman as Brother is to Sister” can be generated through arithmetic operations on the vector representations of these words such that the vector representation of “Brother” - “Man” + “Woman” produces a result which is closest to the vector representation of “Sister” in the model.

Regarding the techniques, **Word2Vec** is one of the most popular strategies to learn word embeddings using a shallow neural network (Mikolov et al. 2013). Additionally, deep learning approaches for contextualized word embeddings such as **BERT** (Devlin et al. 2018) or **EMLO** (Peters et al. 2018) are gaining more popularity. Although there is an emerging trend towards generating embeddings, measuring the quality - i.e. performance - of such embeddings in a different tasks is important and still open problem.

In this project, the Master student will dive deeper into this subject. The study will be addressed using scientific methodology, and standard metrics for the quality evaluating of embeddings in the task of finding semantically related documents.

## TASK

- Literature Review: get to know the state-of-the-art for Language Models.
- Hands-On Machine Learning: Learning of different methods.
- Evaluation: a comparison in different tasks.

## PROJECT PROFILE

- Analysis: 3/5
- Implementation: 5/5 (Python required)
- Literature: 2/5

## HOW TO APPLY

You can apply for this research opportunity by email. Please include your CV, transcripts of UG in your application.

## CONTACT:

Malte Ostendorff | DFKI [malte.ostendorff@dfki.de](mailto:malte.ostendorff@dfki.de), Dr. Vinicius Woloszyn [woloszyn@tu-berlin.de](mailto:woloszyn@tu-berlin.de)