



# Unsupervised Domain adaptation for Sentence Classification

Bernard Kuchler, Jan Šuklje, and Luka Šveigl

## Abstract

Advisors: Slavko Žitnik

## Introduction

We seek to improve article representation in specialized domain articles by adapting pretrained models such as SBERT which are not tuned to specific fields. The non tuned models will serve as a baseline to improve upon. We will fine tune the SBERT model on the SentiNews database [1] for articles related to the keywords: Economy, Economics and Bank. We will attempt to accomplish this by using two adaptation techniques: TSDAE (Transformer-based Denoising AutoEncoder) and GPL (generative pseudo labeling).

SBERT, created by Reimers and Gurevych [2] is an improvement upon BERT for the creation of semantically meaningful sentence embeddings. It will serve as a good baseline as it can encode semantics for sentences making it easier to compare them. Similar sentences have embeddings close together in the vector space while dissimilar sentences are far apart. TSDAE, created by Wang et al. [3] is a domain adaptation and pre-training method. TSDAE is used to introduce noise to input sequences by deleting or swapping tokens and turning them into sentence vectors. It then attempts to reconstruct the original input sequence from the vectors. This allows the model to better capture the characteristics of the target domain GPL, created by Wang et al.[4] is also a unsupervised domain adaptation method. It consists out of query generation, negative mining and pseudo labeling in a data preparation step and a fine-tuning step.

## Methods

## Results

## Discussion

## Acknowledgments

## References

- [1] Boshko Koloski, Matej Martinc, Ilija Tavchioski, Blaž Škrlić, and Senja Pollak. Slovenian keyword extraction dataset from SentiNews 1.0, 2022. Slovenian language resource repository CLARIN.SI.
- [2] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- [3] Kexin Wang, Nils Reimers, and Iryna Gurevych. Tsdae: Using transformer-based sequential denoising auto-encoder for unsupervised sentence embedding learning. *arXiv preprint arXiv:2104.06979*, 2021.
- [4] Kexin Wang, Nandan Thakur, Nils Reimers, and Iryna Gurevych. Gpl: Generative pseudo labeling for unsupervised domain adaptation of dense retrieval. *arXiv preprint arXiv:2112.07577*, 2021.