

Natural Language Processing for the Support of Clinical Depression Detection

Degree programme : BSc in Computer Science | Specialisation : Data Engineering
Thesis advisor : Prof. Dr. Mascha Kurpicz-Briki
Expert : Han Van der Kleij (SBB AG Informatik)

Mental disorders are increasing rapidly in today's society, especially given the current pandemic situation. Social media platforms offer a place for affected people to share their opinions and experiences when dealing with depression. Can Natural Language Processing help clinical practitioners to detect depression?

Introduction

The data on social media platforms provide an important opportunity to develop new technologies that can be valuable tools to professionals in a clinical setup. A completely anonymized dataset of **German** texts from social media containing **specific hashtags** has been **collected and annotated** in preliminary work, which provides the foundation for this thesis.

Goal

The goal of the thesis was to train and fine-tune a machine learning model and validate its accuracy. In addition, a comparison between **traditional machine learning** and **state-of-the-art machine learning** in the **NLP** domain was conducted to highlight the innovation made.

Methodology

A **Logistic Regression** and **LinearSVC** Model was used as the traditional ML method to achieve the goal. On the other hand, a **BERT** model combining the traditional ML methods was used as a **feature extractor** and a stand-alone **fine-tuned** model.

Results

The results showed the importance of a **balanced** dataset and the associated features to prevent **overfitting**. The **downsampling** of the dataset is a

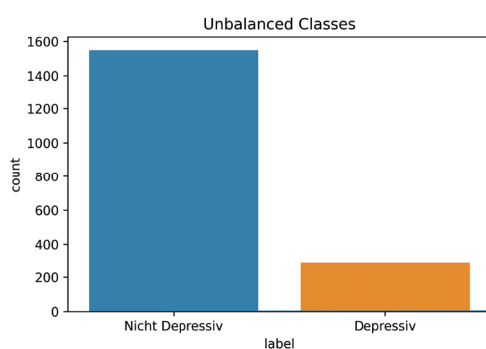
variation to improve overfitting. Moreover, it can be seen that the feature extraction of BERT and the associated **contextual understanding of a language** is an essential milestone in the NLP field.

Conclusion

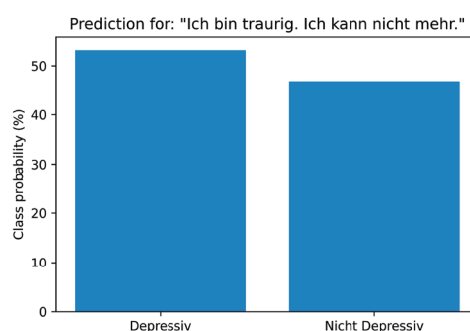
The project has presented many challenges, one of the problems was to find a solution to overfitting without collecting new data. Another challenge was that data from the „real world“ can be messy, especially from social media platforms, which must be kept in mind during **pre-processing**. In addition, the interrater reliability showed that the disagreement between annotators was recognizable in the machine learning model itself.



Kevin Amalathas



Class Distribution



Prediction Example of Fine-tuned BERT with Downsampled Dataset