



Offensive language exploratory analysis

Tomaž Martinčič, Dimitar Stefanov and Žiga Trojer

Abstract

Offensive speech is often present on the Internet. Many times this is not desirable on the pages, or it is forbidden, so we want to detect it automatically. This task has attracted a lot of attention recently, and as a result, many language processing techniques and algorithms have been developed. In this report, we present the capabilities of a number of them. Furthermore, as our end goal we try to produce a meaningful visualization of the distances among the different classes, or better say types of hate speech.

Keywords

hate speech, offensive language, exploratory analysis

Advisors: prof. dr. Slavko Žitnik

Introduction

Our first task was to find as many offensive speech data sets as possible. Out of the 65 datasets for which we were given descriptions and references, we decided to focus on the 25 English language datasets. We had a lot of problems with the protection of personal data - some data sets were not accessible. Many data sets consisted of Twitter messages, but only tweet IDs were given - scrape data is required to analyze these tweets.

Another problem that we encountered was that we found some data sets and most of them had only binary class - hate speech or not. We also found some others, but we can not merge those data sets, because the classes are different. In addition, the datasets were obtained from different media, so formatting issues need to be taken into account in some cases as well.

We have documented our whole analysis of the quality and suitability of the datasets in the Github repository¹. Additionally, we have come across one useful dataset on Kaggle (Twitter), but it has not been added to Github due to its size (6 GB). This dataset is rather abundant, so we used it for one of our initial explorations.

We also found some methods that we could use. We will list some of them, so the decision which one to choose will be easier.

1. Non-contextual pre-trained word embeddings: Word2Vec, Glove, fastText

2. Contextual pre-trained word embeddings: BERT, ELMo
3. Support Vector Machines (SVM) over Bag-of-Words vectors
4. Recurrent Neural Networks with Long Short-term Memories [1]
5. Convolutional Neural Networks (CNN) [2]
6. Logistic regression (LR), and multi-layer perceptrons (MLP) - using bag-of-words representations based on either word or character-level n-grams [3]
7. Lexical Syntactic Feature (LSF), support vector machine with different kernels: Polynomial, Dot, Radial, ANOVA and Epachnenikov. [4]
8. Ensemble models (combination of context-aware logistic regression and context-aware neural networks) [5]

Our plan is that we start with some simple approaches such as TF-IDF (Term Frequency - Inverse Document Frequency) and BoW (Bag of Words). Then we will try to implement some of the other approaches suggested in the project description (Word2Vec, BERT,...), that best fit our problem.

For the visualization part, we plan on performing PCA (Principal Component Analysis), MDS (Multi-Dimensional Scaling) or representation of the embeddings with t-SNE.

¹https://github.com/tm1897/nlp_offensive_lang

