

Heidelberg University
Institute of Computer Science
Database Systems Research Group

Project Proposal for the lecture Text Analytics
Hate Speech Detection

Team Member: Christoper Klammt, 3588474, Applied Computer Science
iv249@stud.uni-heidelberg.de

Team Member: Felix Hausberger, 3661293, Applied Computer Science
eb260@stud.uni-heidelberg.de

Team Member: Nils Krehl, 3664130, Applied Computer Science
pu268@stud.uni-heidelberg.de

1 Motivation

Legal implementations on handling hate speech is different from one country to another. While hate speech is not prohibited in the United States due to the *freedom of speech*, other countries especially in the European Union can sue hate speech actors for either offending the public order or human dignity. While being able to prosecute actors in public without much effort, the internet and especially social media platforms provide an easy and anonymous way to practice hate speech without legal consequence enforcements. Several steps were taken to tackle hate speech online, one of them being the *code of conduct* on countering illegal hate speech online, an initiative of the European Commission in close collaboration with major IT companies like *Facebook*, *Microsoft*, *Twitter* and *YouTube* [1]. While respecting the freedom of speech, these companies commit to delete hate speech contributions within 24 hours of the initial deletion request. To further automatize the process of detecting hate speech contributions, several text analytics approaches have been evaluated in the recent past. Many of them are using methods of *natural language processing* and *deep learning* for hate speech detection and rely on meaningful features being learned automatically by neural networks instead of using hand-crafted features. In this work, the boundaries of conventional text analytics approaches for hate speech detection including manual feature selection and subsequent text classification of hate speech and non hate speech documents for the Twitter API should get evaluated. The result of this work should show which features work best for which classifier and which problems can be addressed with conventional text analytics methods and which not.

To add: -papers -datasets

2 Research Topic Summary

There are three aspects to consider, derived from the research on automatic hate speech detection:

1. Raw dataset
2. Feature extraction
3. Machine Learning technique

The following sections introduce these three aspects.

2.1 Raw dataset

To obtain a dataset one can use an existing already labelled dataset (e.g. done by Watanabe, Bouazizi and Ohtsuki [2]), or one could label the data by hand (e.g. made by Oriola and Kotzé [3]). Depending on the dataset the classes in which the data is classified can differ. Prominent classifications for datasets in the domain of hate speech are the binary classification (no hate speech, hate speech) and the ternary classification (clean, offensive, hate speech). TODO: Beispieldatensätze

Another aspect is the distribution between the different classes. Is the distribution equally among the different classes or not? In case of imbalanced distributions the machine learning approach can be less performant and accurate [3].

2.2 Feature extraction

The main difference between the two introduced methodologies (neuronal network, classical Machine learning techniques) is the process of feature extraction. In neuronal network approaches the used features are learned automatically, whereas classical Machine Learning techniques require a manual feature extraction process. This is highly based on text analytics.

The following list provides an overview of possible features and their technical methods from the area of text analytics. The list is clustered into the four categories based on Watanabe, Bouazizi and Ohtsuki [2]. The possible text analytics approaches derive from the literature review of Fortuna and Nunes [4].

- **Sentiment-based features:** Is the tweet rather positive or negative?
Text analytics approaches: Dictionaries, Rule Based Approaches, Objectivity-Subjectivity of the Language, Declarations of Superiority of the Ingroup
- **Semantic features:** Which parts of the tweet are emphasized?
Text analytics approaches: TF-IDF, Part-of-speech, Profanity Windows, Lexical Syntactic Feature-based, Topic Classification, Template Based Strategy, Word Sense Disambiguation Techniques, Othering Language
- **Unigram features:** Are there any specific words marking hate speech?
Text analytics approaches: N-grams, Bag-of-words
- **Pattern features:** Are there any specific patterns marking hate speech?
Text analytics approaches: Part-of-speech, Dictionaries, Typed Dependencies, Word Embeddings

2.3 Machine Learning technique

Current research in the area of hate speech detection is often built upon neuronal network techniques. Exemplary therefore are the works of Roy et al. [5], Setyadi, Nasrun and Setianingsih [6] and Kapil, Ekbal and Das [7], to name just a few examples. Among other metrics by considering the accuracy one can evaluate the success of a chosen approach. Roy et al. [5] have reached accuracies up to 95% for detecting whether a tweet is hateful or not.

Nevertheless even classical Machine Learning techniques such as Support Vector Machines (SVM) combined with text analytics methods are a promising approach. Watanabe, Bouazizi and Ohtsuki [2] used a decision tree and reached an accuracy of 87.4%.

Furthermore the results of classical Machine Learning techniques can be used as reference value to evaluate the neuronal network approach. This was done by Roy et al. [5].

3 Project Description

References

- [1] European Commission, *The code of conduct on countering illegal hate speech online*, European Commission Homepage, 2020. [Online]. Available: https://ec.europa.eu/commission/presscorner/detail/en/qanda_20_1135 (visited on 11/19/2020) (cit. on p. 1).
- [2] H. Watanabe, M. Bouazizi, and T. Ohtsuki, "Hate speech on twitter: A pragmatic approach to collect hateful and offensive expressions and perform hate speech detection," *IEEE Access*, vol. 6, pp. 13 825–13 835, 2018. DOI: 10.1109/ACCESS.2018.2806394 (cit. on pp. 2, 3).
- [3] O. Oriola and E. Kotzé, "Evaluating machine learning techniques for detecting offensive and hate speech in south african tweets," *IEEE Access*, vol. 8, pp. 21 496–21 509, 2020. DOI: 10.1109/ACCESS.2020.2968173 (cit. on p. 2).
- [4] P. Fortuna and S. Nunes, "A survey on automatic detection of hate speech in text," *ACM Comput. Surv.*, vol. 51, no. 4, Jul. 2018, ISSN: 0360-0300. DOI: 10.1145/3232676. [Online]. Available: <https://doi.org/10.1145/3232676> (cit. on p. 2).

- [5] P. K. Roy, A. K. Tripathy, T. K. Das, and X. -Z. Gao, “A framework for hate speech detection using deep convolutional neural network,” *IEEE Access*, vol. 8, pp. 204 951–204 962, 2020. DOI: 10.1109/ACCESS.2020.3037073 (cit. on p. 3).
- [6] N. A. Setyadi, M. Nasrun, and C. Setianingsih, “Text analysis for hate speech detection using backpropagation neural network,” in *2018 International Conference on Control, Electronics, Renewable Energy and Communications (ICCEREC)*, 2018, pp. 159–165. DOI: 10.1109/ICCEREC.2018.8712109 (cit. on p. 3).
- [7] P. Kapil, A. Ekbal, and D. Das, *Investigating deep learning approaches for hate speech detection in social media*, 2020. arXiv: 2005.14690 [cs.CL] (cit. on p. 3).