

Paper Title:

DEEP LEARNING FOR HATE SPEECH DETECTION: A COMPARATIVE STUDY

Paper Link:

<https://arxiv.org/pdf/2202.09517.pdf>

1. Summary**1.1 Motivation**

The motivation behind this report stems from the prevalent challenge of identifying and categorizing offensive language within the realm of social media. Its core purpose is to build upon existing methodologies and frameworks, aiming to augment the accuracy of detecting offensive content. By leveraging prior research findings, the report endeavors to contribute substantially to this domain.

1.2 Contribution

The primary contribution of this endeavor lies in achieving a commendable accuracy rate of 74% in the classification of offensive tweets. This accomplishment is coupled with a critical evaluation of impending challenges within the field. It identifies and delineates crucial areas that demand further attention and enhancement.

1.3 Methodology

The methodology adopted in this pursuit integrates insights gleaned from prior research. It involves meticulous data balancing, refined feature selection tailored for classification, and a series of preprocessing stages aimed at fine-tuning the overall approach. This methodical process has been instrumental in achieving the reported accuracy and in providing insights into the challenges faced.

1.4 Conclusion

Concluding this study acknowledges and celebrates the successes achieved in offensive tweet classification. However, it also underscores persistent obstacles and limitations. It emphasizes the need for standardized definitions and calls for concerted efforts to address the ongoing challenges for future advancements in this field.

2. Limitations**2.1 First Limitation**

- Short length and lack of context in tweets impede accurate classification.
- Removal of identifiers reduces precision in identifying offensive content.

2.2 Second Limitation

- Difficulty in detecting sarcasm and determining offensive intent.
- Lack of standardized definitions and data inconsistencies pose challenges.

3. Synthesis

This report's synthesis reflects on its broader implications and potential applications. The envisioned future scope involves the development of sophisticated classifiers capable of

accommodating diverse multimedia content for more holistic detection. Overcoming the challenges posed by multilingual communication nuances and privacy concerns will necessitate the evolution of nuanced, adaptive algorithms.

The potential applications extend to leveraging advanced techniques like Word2Vec and exploring representation learning to augment the accuracy of classification. This proactive approach stands poised to revolutionize the landscape of offensive content detection on social media platforms.