# Offensive Text Detection Using Classical Ai/Ml Techniques

**Nandini Tidke[1], Komal Chopade[2], Gaurav Kalsait[3], Jai Singh[4], Prof. Santosh Biradar[5]**
Students, Department of Computer Engineering[1, 2,3,4]
Assistant Professor, Department of Computer Engineering[5]
D. Y. Patil College of Engineering, Pune, Maharashtra, India

**Abstract:** *Due to the substantial growth of internet users and its spontaneous access via electronic devices, the amount of electronic contents has been growing enormously in recent years through instant messaging, social networking posts, blogs, online portals and other digital platforms. Unfortunately, the misapplication of technologies has increased with this rapid growth of online content, which leads to the rise in suspicious activities. People misuse the web media to disseminate malicious activity, perform the illegal movement, abuse other people, and publicize suspicious contents on the web. The suspicious contents usually available in the form of text, audio, or video, whereas text contents have been used in most of the cases to perform suspicious activities. Thus, one of the most challenging issues for NLP researchers is to develop a system that can identify suspicious text efficiently from the specific contents. We define this task as being able to classify a tweet as offensive or not. A set of ML classifiers with various features has been used on our developed corpus, consisting of 7000 English text documents where 5600 documents used for training and 1400 documents used for testing. The performance of the proposed system is compared with the human baseline and existing ML techniques.*

**Keywords:** Test detection

## REFERENCES

[1]. W. Warner and J. Hirschberg. (2012). Detecting hate speech on the world wide web.Proceeding LSM '12 Proc. Second Work. Lang. Soc.Media, no. Lsm, pp. 19–26.

[2]. I. Kwok and Y. Wang.(2013). Locate the hate: detecting tweets against blacks.Twenty-Seventh AAAI Conf. Artif. Intell., pp. 1621–1622.

[3]. I. Alfina, D. Sigmawaty, F. Nurhidayati, and A. N. Hidayanto.(2017). Utilizing hashtags for sentiment analysis of tweets in the political domain. In Proceedings of the 9th International Conference on Machine Learning and Computing, pp. 43–47.

[4]. Freund, Y; Schapire, R.E.(1999). Large margin classification using the perceptron algorithm. Machine Learning, 37(3):277–296.

[5]. Kim, Y.H. et al. (2000). Text filtering by boosting naive Bayes classifiers. ACM SIGIR Conference:p168-175.

[6]. Parikh R, Movassate M. (2009). Sentiment analysis of user-generated Twitter updates using various classification techniques. CS224N Final Report;pages. 1–18.

[7]. Pak A, Paroubek P. (2010). Twitter as a Corpus for Sentiment Analysis and Opinion Mining. In: LREC. vol. 10; pages. 1320–1326.

[8]. Gaudette L, Japkowicz N. (2009). Evaluation methods for ordinal classification.In Advances in Artificial Intelligence. Springer; p. 207–210.

[9]. Go A, Bhayani R, Huang L. (2009). Twitter sentiment classification using distant supervision. CS224N Project Report, Stanford;p. 1–12