

# Project Proposal Group 14

Joost Bambacht (4025016), Robin Faber (4560175), Wessel Turk (4599004), Wouter Zorgdrager (4472977)

## 1 Problem description

The problem we will try to tackle during this project is clickbait detection. This problem is interesting due to the increasing role the Internet plays in our daily lives. Due to the enormous presence of advertisements online nowadays, websites try their best to get people to visit their pages. This can lead to titles that are meant to engage users, sometimes purposefully written in a misleading way to attract users. Obviously it is important to be able to detect and prevent clickbait in order to ensure that users get the information they are expecting to get.

The paper that we will reproduce is the one given on the course page for the clickbait challenge of 2017 [1]. The approach of this paper is to extract features from a (labeled) dataset on clickbait articles. These features are then used to train a classifier to determine whether an article is clickbait or not.

## 2 Resources

In the original paper [1], they made use of two datasets. One smaller initial training set and one bigger set which, besides training, can be used to validate and test the model. We will use these same datasets. To extract features from these sets, we will write our own Python scripts. Finally, to train the ML model (classifier) we will use of on the popular machine learning libraries like scikit <sup>1</sup>.

## 3 Methodology

The main goal of this project is to reproduce Table 2 of the original paper [1]. This gives for each machine learning algorithm the various performances of the evaluation metrics on the training set and validation set. More about the evaluation metrics themselves is explained in Section 4. In the original paper, the information gain of each feature is calculated and shown in Table 1. We believe this is most likely out of the scope of this project and will thus not be reproduced.

## 4 Evaluation

To make sure the results in our evaluation are legitimate we will split the dataset in a training and test set. We will only use the test set in the end to determine the performance of our model without incorporating this set in the training phase. To evaluate the performance of this test set we will use the following performance measures: AUC, precision, recall and accuracy. These are the measures used in our chosen paper [1]. Finally, we will also evaluate using the F1 score to express the balance between the recall and precision measure.

## 5 Background reading

Paper [1] is the paper that will be reproduced in this project. The other papers in the references were found to be relevant to this study and will (possibly) be used as background information.

---

<sup>1</sup><https://scikit-learn.org/>

## References

- [1] A. Elyashar, J. Bendahan, and R. Puzis, “Detecting clickbait in online social media: You won’t believe how we did it,” 2017.
- [2] D. Daoud and S. El-Seoud, “An effective approach for clickbait detection based on supervised machine learning technique,” *International Journal of Online and Biomedical Engineering (iJOE)*, vol. 15, p. 21, February 2019.
- [3] A. Chakraborty, B. Paranjape, S. Kakarla, and N. Ganguly, “Stop clickbait: Detecting and preventing clickbaits in online news media,” October 2016.
- [4] A. Chakraborty, R. Sarkar, A. Mrigen, and N. Ganguly, “Tabloids in the era of social media? understanding the production and consumption of clickbaits in twitter,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 1, no. CSCW, pp. 1–21, 2017.
- [5] P. Biyani, K. Tsioutsoulis, and J. Blackmer, ““8 amazing secrets for getting more clicks”: Detecting clickbaits in news streams using article informality,” in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [6] A. Pujahari and D. S. Sisodia, “Clickbait detection using multiple categorisation techniques,” *Journal of Information Science*, 2019.
- [7] N. Zuhroh and N. Rakhmawati, “Clickbait detection: A literature review of the methods used,” *Register: Jurnal Ilmiah Teknologi Sistem Informasi*, vol. 6, p. 1, 10 2019.