

Comparitive analysis of clickbait classifiers

Pradyumna Heddur Nagendra, Priya Chaudhary

Section 1 Introduction

Clickbait is a term that describes deceiving web content that uses ambiguity to provoke the user into clicking a link. These baits may trick the readers into clicking, but in the long run, clickbait usually doesn't live up to the expectation of the readers and leave them disappointed.



Section 3 Dataset

Stop Clickbait Dataset is used in this project. It contains 16,000 article headlines categorized as "clickbait" and "non-clickbait". The clickbait articles have been pulled from websites including Buzzfeed and Upworthy, while the non-clickbait articles come from sites including Wikinews, The New York Times, and The Guardian.

	headline	clickbait
0	Should I Get Bings	1
1	Which TV Female Friend Group Do You Belong In	1
2	The New "Star Wars: The Force Awakens" Trailer...	1
3	This Vine Of New York On "Celebrity Big Brothe...	1
4	A Couple Did A Stunning Photo Shoot With Their...	1
...
31995	To Make Female Hearts Flutter in Iraq, Throw a...	0
31996	British Liberal Democrat Patsy Calton, 56, die...	0
31997	Drone smartphone app to help heart attack vict...	0
31998	Netanyahu Urges Pope Benedict, in Israel, to D...	0
31999	Computer Makers Prepare to Stake Bigger Claim ...	0

32000 rows x 2 columns

Section 5 Evaluation

The different classification models are evaluated using the accuracy, recall, F1-score, and confusion matrix.

TF-IDF:

Method	Accuracy	Recall	F1-Score
LSTM	0.961	0.953	0.961
SVM	0.962	0.962	0.962
Random Forest	0.912	0.969	0.917

Word2Vec:

Method	Accuracy	Recall	F1-Score
LSTM	0.959	0.952	0.958
SVM	0.812	0.720	0.797
Random Forest	0.825	0.770	0.816

BERT:

Accuracy	Recall	F1-Score
0.828	0.782	0.819

References

[1] Abhijnan Chakraborty, Bhargavi Paranjape, Sourya Kakarla, and Niloy Ganguly. Stop clickbait: Detecting and preventing clickbaits in online news media. *CoRR*, abs/1610.09786, 2016.

[2] Mark Bronakowski, Mahmood Al-khassaweneh, and Ali Al Bataineh. Automatic detection of clickbait headlines using semantic analysis and machine learning techniques. *Applied Sciences*, 13(4), 2023.

[3] Suhaib R. Khater, Oraib H. Al-sahlee, Daoud M. Daoud, and M. Samir Abou El-Seoud. Clickbait detection. In *Proceedings of the 7th International Conference on Software and Information Engineering*, ICSIE '18, page 111–115, New York, NY, USA, 2018. Association for Computing Machinery.

Acknowledgement

We would like to convey our sincere gratitude to Professor Paul Plöger and Tim Metzler for giving us the chance to collaborate on this project. It helped us do in-depth research and increase our understanding of NLP.

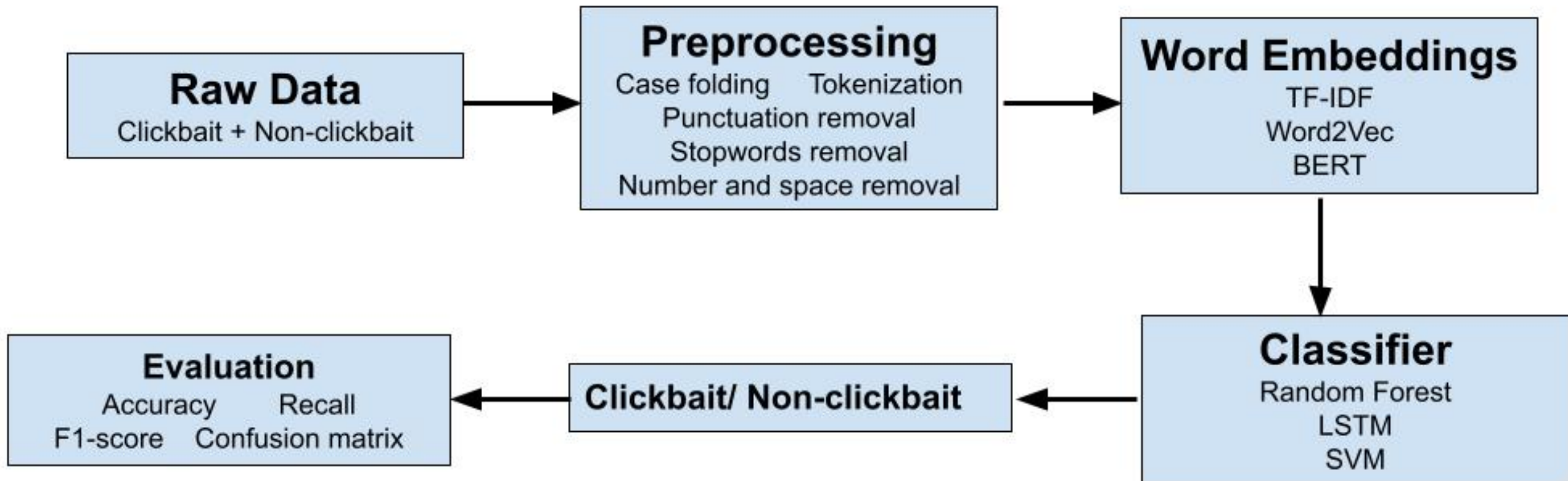
Section 2 Methodology

To prepare the raw data for subsequent processing, clickbait and non-clickbait headlines are first processed. Tokenization, stop word removal, lower-case processing, and punctuation removal are some of the processing techniques.

For comparative analysis, three classification models (SVM, LSTM, and Random forest) are combined with embedding techniques like TF-IDF and Word2Vec.

For Classification, the BERT model is also employed.

Section 4 Network Architecture



Contact

Pradyumna Heddur Nagendra, Priya Chaudhary
Hochschule Bonn-Rhein-Sieg
Email:pradyumnaheddurnagendra@smail.inf.h-brs.de,
priya.chaudhary@smail.inf.h-brs.de