



Fake News Detection Using Machine Learning

Prof. Supriya Yadav¹, Sachin Tendulkar², Aditya Sakore³, Shrushti Umalkar⁴, Shivam Kobal⁵

^{1,2,3,4,5}Department of Information Technology, SPPU

Abstract— Malicious URLs have been widely used to mount various cyber-attacks including spamming, phishing and malware. Detection of some unreal URLs and identify their threat types are critical to handle. Existing methods typically detect URLs of only one attack type. Our method uses a variety of link structures, articles, journals, web page contents, DNS information, and network traffic. Our experimental studies with so many benign URLs and few malicious URLs obtained from real-life Internet sources show that our method delivers a superior performance.

Keywords— URL, Web Threats, Cyber-attacks, Spamming

I. INTRODUCTION

While the World Wide Web has become a killer application on the Internet, it has also brought in an immense risk of cyber-attacks. Adversaries have used the Web as a vehicle to deliver malicious attacks such as phishing, spamming, and malware infection. For example, phishing typically involves sending an email seemingly from a trustworthy source to trick people to click a URL (Uniform Resource Locator) contained in the email that links to a counterfeit web page. A common countermeasure is to use a blacklist of malicious URLs, which can be constructed from various sources, this work was done when Hyunsang Choi was an intern at Microsoft Research Asia. Contact author: Bin B. Zhu (binzhu@ieee.org) particularly human feedbacks that are highly accurate yet time-consuming. Blacklisting incurs no false positives, yet is effective only for known malicious URLs. It cannot detect unknown malicious URLs. The very nature of exact match in blacklisting renders it easy to be evaded.

This weakness of blacklisting has been addressed by anomaly-based detection methods designed to detect unknown malicious URLs. In these methods, a classification model based on discriminative rules or features is built with either knowledge a priori or through machine learning. Selection of discriminative rules or features plays a critical role for the performance of a detector. A main research effort in malicious URL detection has focused on selecting highly effective discriminative features. Existing methods were designed to detect malicious URLs of a single attack type, such as spamming, phishing, or malware.

II. OUR FRAMEWORK

Our method consists of three stages: training data collection, supervised learning with the training data, and malicious URL detection and attack type identification. These stages can operate sequentially as in batch learning, or in an interleaving manner: additional data is collected to incrementally train the classification models while the models are used in detection and identification. Interleaving operations enable our method to adapt and improve continuously with new data, especially with online learning where the output of our method is subsequently labelled and used to train the classification models.

Learning Algorithms

The two tasks performed by our method, detecting malicious URLs and identifying attack types, need different machine learning methods. The Logistic Regression is used to detect malicious URLs. The second task is a multi-label classification problem. Two multi-label classification methods are used to identify attack types.

Automated Fake News Detection

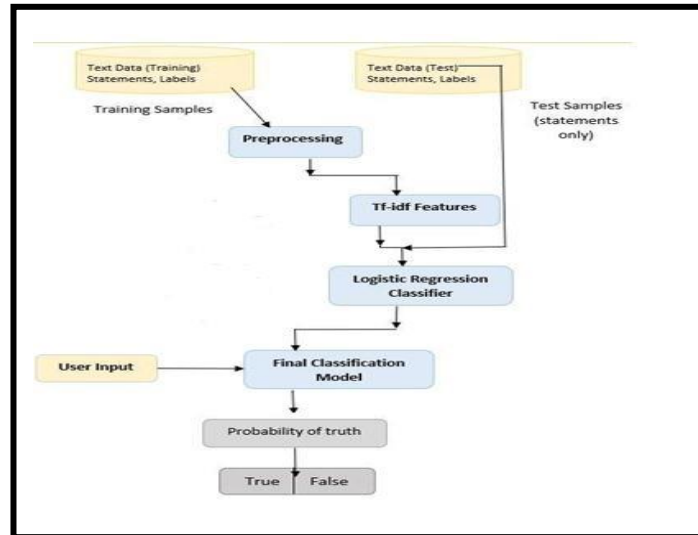
One of the most obvious applications of our datasets is to facilitate the development of machine learning models for automatic fake news detection. In this task, we frame so many way multiclass text classification problem. And the research questions are: 1] Based on surface-level linguistic realizations only, how well can machine learning algorithms classify a short statement into a fine-grained category of fakeness?

III. SYSTEM ARCHITECTURE

A. Steps to initialize

We randomly initialize a matrix of embedding vectors to encode the metadata embedding in our project. We use a convolution layer to capture the dependency between the meta-data vector(s). We then concatenate the max-pooled text

representations with the meta-data representation from the bi-directional LSTM, and feed them to fully connected layer with a activation function to generate the final prediction.



We used pre trained so many Google News to start the text embedding. We put all the hyperparameters to the validation set

B. Next Stage while processing the data

The evaluation results on the LIAR dataset. The top section: text-only models. The bottom: text + meta-data hybrid models. while no L2 penalty was imposed. We considered 0.5 and 0.8 as dropout probabilities. The hybrid model requires 5 training epochs. . . First, we compare models using so many type of text features. We see that the majority baseline on this dataset gives about 0.204 and 0.208 accuracy on the validation and test sets respectively. Due to overfitting, that algorithm not performs well. The CNNs outperformed all models, resulting in an accuracy of 0.270 on the holdout test set. We compare the predictions from the CNN model with logistic regression via a two-tailed paired t-test, and logistic regression was significantly better .

Machine learning approach

Detection of single attack type. Machine learning has been used in several approaches to classify malicious URLs. They used site dependent heuristics, such as words used in a page or title and fraction of visible uses URLs in emails as input and outputs regular expression signatures that can detect spam .They used a large publicly available corpus of legitimate and phishing emails. Their classifiers examine ten different features such as the number of URLs in an e-mail, the number of domains and the number of dots in these URLs

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

subject to

$$\sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C, i = 1, 2, \dots, n$$

Existing machine learning-based approaches usually focus on a single type of malicious behaviour. They all use machine learning to tune their classification models. Our method is also based on machine learning, but a new and more powerful and capable classification model is used. In addition, our method can identify attack types of malicious URLs. These innovations contribute to the superior performance and capability of our method.

**IV. ALGORITHM**

Algorithm for Fake News Detection :

- Data Preprocessing stage
- Data training and testing stage
- In this stage we try to use so many different machine learning algorithms like svm,,random forest,logistic
- Train the model finally on logistic regression to get the result.
- Prepare the prediction that is if news is fake or true on the basis of logistic regression
- Stop.

V. CONCLUSION AND FUTURE SCOPE

We introduced LIAR, a new dataset for automatic fake news detection. Compared to prior datasets, LIAR is an order of a magnitude larger, which enables the development of statistical and computational approaches to fake news detection. LIAR's authentic, real-world short statements from various contexts with diverse speakers also make the research on developing broad-coverage fake news detector possible. We show that when we have to combine the meta-data with text, so much of improvements can be achieved for this fake news detection. It is also possible to explore the task of automatic fact-checking over knowledge base in the future. Our corpus can also be used for stance classification, argument mining, topic modelling, rumour detection, and political NLP research. In future scope is use to live detection of news , if any scam or fake news is appeared on the social media then we directly send the image of that post to the server or near police station or cyber cell . To take action against such people.

REFERENCES

- [1] Ronan Collobert, Jason Weston, Leon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Shubhendu Apoorv, Sudharshan Kumar Bhowmick, R Sakthi Prabha, "Indian sign language interpreter using image processing and machine learning". 2011. Natural language processing (almost) from scratch. *Journal of Machine Learning Research* 12(Aug):2493–2537.
- [2] Koby Crammer and Yoram Singer. 2001. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of machine learning research* 2(Dec):265–292.
- [3] Song Feng, Ritwik Banerjee, and Yejin Choi. 2012. Syntactic stylometry for deception detection. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers-Volume 2*. Association for Computational Linguistics, pages 171–175.
- [4] William Ferreira and Andreas Vlachos. 2016. Emergent: a novel data-set for stance classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. ACL.
- [5] Alex Graves and Jurgen Schmidhuber. 2005. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks* 18(5):602–610.