



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 9 Issue: IV Month of publication: April 2021

DOI: <https://doi.org/10.22214/ijraset.2021.33900>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Fake News Detection using Machine Learning

Abha Tewari¹, Sujoy Mitra², Ishaan Nangrani³, Pratik Nathani⁴, Alish Wadhvani⁵

^{1, 2, 3, 4, 5}Department of Computer Engineering, Vivekanand Education Society's Institute of Technology

Abstract: *The ubiquity of fake news today is spreading like an invisible forest fire which means we can't distinguish between the real and fake news so easily. People believe in anything the social media, news websites, online newspapers, Blog/Vlog/Weblog posts, etc. show very easily. Therefore the credibility, integrity and authentication of news are imperative in every field. Data or information, today is travelling faster than the speed of light so we need to be quick to assess it as well. This has made certain computational, logical and analytical tools that can help us identify real news from the online content. In this paper, we have used a dataset to recognise false news. Pre-processing, feature extraction, classification, and prediction are all elucidated in detail. Some operations such as tokenizing, stemming, and Data Exploration such as response variable distribution and data quality check are performed by the pre-processing functions (i.e. null or missing values). Function extraction techniques include quick bag-of-words, n-grams, and TF-IDF. For fake news detection with a probability of truth, a logistic regression model is used as a classifier.*

Keywords: *Fake news detection, Logistic regression, TF-IDF vectorization, Tokenization, Stemming, Lemmatization, n-grams, Bag-of-Words, Probability of Truth*

I. INTRODUCTION

Modern life is becoming more and more reliant on the internet and cutting edge technological advancements to deliver the latest and the most trending/viral information to each and every individual around the globe but "Every coin has two sides". As the development in information transfer continues to advance towards a better future, it also comes with the growth in the pessimistic side for example Trolling, Fake News Spreading, Spamming, etc. which has to stop or at least mitigate as much as possible. Facebook, Twitter, Reddit, Youtube, and Whatsapp are some of the social media sites that are used to apportion fake news [1].

The Accuracy of the model is at the most 70-85% on most of the models. We have included the Naive Bayes classifier, Linguistic features based, Bounded decision tree model, SVM, etc. The objective of this paper is to ameliorate the accuracy of recognising fake news beyond what is already present. By fabricating a new model which will conclude the spurious news articles on the basis of the following criteria: spelling mistake, jumbled sentences, punctuation errors etc.

II. LITERATURE SURVEY

There are three categories of fake news in general. The first one is the fake news, which is not backed by any kind of research and completely made up by the authors just to make it attractive. The second category is fake satire news, which is fake news with a primary goal of amusing people. The third category is the poorly written news articles, which are just partially true and although they contain some real news, but are not completely true.

Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu explored in their article the two phases of reviewing the fake news: characterization and detection. The basic concepts of fake news were in the characterization phase and the detection phase had a review of existing fake news detection approaches. These include feature extraction and construction of model [2].

In the paper by Hadeer Ahmed, Issa Traore, and Sherif Saad, a comparison was made between different feature extraction techniques and machine classifiers. Their model used n-gram analysis and machine learning techniques. They obtained a maximum accuracy of 92% using the Term Frequency-Inverse Document Frequency (TF-IDF) and Linear Support Vector Machine (LSVM) [3].

In order to identify fake contents in online news, Perez-Rosas, Veronica & Kleinberg, Bennett and Lefevre Alexandra and Rada Mihalcea had used two different datasets. Different classification models were developed by them in order to get the maximum accuracy. They used linear sum classifier and fivefold cross verification in order to get the accuracy, precision and recall and FI scores averaged over the five iterations [4].

E.M Okoro, B.A Abara, A.O. Umagba, A.A. Ajonye and Z. S. Isa in their publication had combined human-based and machine-based approaches since these both cannot solve the problem of human literacy on their own. They introduced a Machine Human (MH) model to detect fake news in social media. This model combines the human literacy news detection tool and machine linguistic and network-based approaches. So these two approaches worked simultaneously with each other to detect fake news [5].

III.EXISTING SYSTEM AND THEIR LACUNA

There are various methods to resolve the issue of fake news. We've approached it by using various NLP techniques. The main cause behind this fake news are complex issues. Performing automatic fact checking with the help of classification techniques can also be done. The problem lies at classification of legitimate and fake instances. The classification of text mainly depends on the linguistic characteristics of the full text. This can be done with the help of various machine learning algorithms as well as models which depend on word vectors.

The feature based classification techniques involved are showing significant results with the help of supervised models. Features based linguistically from the text are extracted on the basis of levels i.e. characters, words and sentences. The problem lies in finding the fake articles. It may be difficult to tag an article with a fake remark. It can also be possible that the article is being written by an amateur. Also, it might be possible that it is not written by a journalist. At times there are instances when professional journalists tamper with the article for personal gains [6].

It is really difficult to define fake news. In simple words, an article containing wrong information can be called a fake news and an article with verified information can be tagged as true. But when we dig deeper into this matter we came to know that we cannot just call a news article fake or real. It might be possible that the article is partly fake or it may be partly real. Various features to be taken into account during fake news detection are as follows:

- 1) *Source*: The publisher of the article.
- 2) *Headline*: A headline can at times play a major role in the classification.
- 3) *Body Text*: This is the most important part which needs to be analysed in order to get the desired results.

IV.PROPOSED ARCHITECTURE

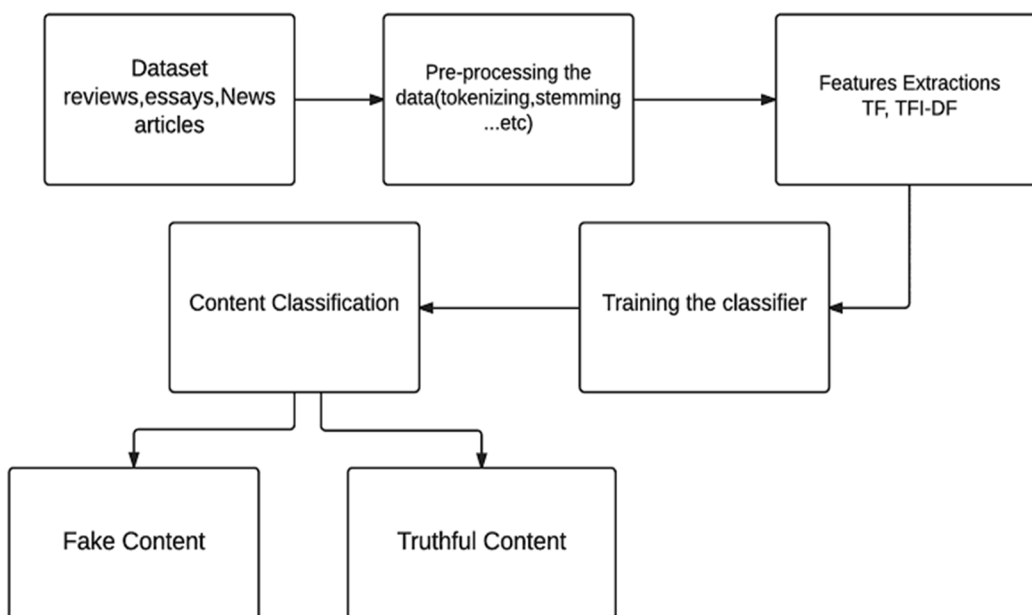


Fig. 1 Proposed Architecture

A. Data Pre-Processing

This contains all the pre-processing functions needed to refine and process all the input data and texts. First the train, test and validation data files are read and then some preprocessing tasks like tokenizing, stemming etc are performed. Then some data quality checks are done like null or missing values etc. and some exploratory data analysis is performed like response variable distribution.

B. Stemming

Stemming is the method of reducing a word to its word stem that affixes to suffixes and prefixes or to the roots of words referred to as lemma. Stemming is very important in natural language processing (NLP). The stem needn't be a clone of the morphological root of the word.

C. Tokenization

Tokenization is essentially dividing a phrase, sentence, paragraph, or a whole text document into smaller units, like individual words or terms. Each of these smaller units are called tokens. These tokens facilitate in understanding the context or developing the model for the NLP. The tokenization helps in decoding the meaning of the text by analyzing the sequence of the words.

D. Feature Selection

In this section we've performed feature extraction and selected methods from sci-kit learn python libraries. For feature selection, we've used strategies like easy bag-of-words and n-grams and then use term frequency like TF-IDF weighting.

E. TF-IDF Vectorizer

TF-IDF (Term Frequency-Inverse Document Frequency) Vectorizer is an analytical measure that gauges how pertinent a word is to a document in a collection of documents. This is done by multiplying two metrics: how many times a word emerges in a document, and the inverse document frequency of the word across a set of documents. It has many applications, essentially in automated text analysis, and is advantageous for scoring words in machine learning algorithms for Natural Language Processing (NLP).

F. Multinomial Naive Bayes Algorithm

Multinomial Naive Bayes algorithm is a probabilistic learning methodology that's principally utilized in Natural Language Processing (NLP). The algorithmic program relies on the Bayes theorem and predicts the tag of a text like a bit of email or newspaper article. It calculates the probability of each tag for a given sample and then provides the tag with the highest probability as output. Naive Bayes classifier is an assortment of many algorithms where all the algorithms share one common principle, which is that every feature being classified isn't associated with any other feature. The presence or absence of a feature doesn't affect the presence or absence of the other feature [8].

V. RESULTS

We have used Multinomial Naive Bayes classifier which will serve the model and work with the user input. Here, we have presented a detection model for fake news using TF-IDF analysis. We have investigated different machine learning techniques. The proposed model achieves accuracy of approximately 85% when using TF-IDF vectorizer and Multinomial Naive Bayes classifier.

Table I
Classification Report

	Precision	Recall	F1 Score	Support
Fake	0.96	0.73	0.83	618
Real	0.79	0.97	0.87	649
Accuracy			0.85	1267
Macro Average	0.87	0.85	0.85	1267
Weighted Average	0.87	0.85	0.85	1267

VI. CONCLUSION

The pre-processing, training and prediction phases have been implemented and the performance of the model has been assessed using accuracy, precision and recall. The model fulfills its task of detecting fake news.

VII. FUTURE SCOPE

This model can be more applicable in the future for other regional languages (Like Hindi, Marathi, Bengali, etc.) and especially using a native country dataset by either collecting data using a Twitter/Google/Reddit API or by web scraping from a Social Media/News Website (We have already attempted in doing so using an American Dataset found in Kaggle).



REFERENCES

- [1] Shivam B. Parikh and Pradeep K. Atrey, "Media-Rich Fake News Detection: A Survey", IEEE Conference on Multimedia Information Processing and Retrieval, 2018.
- [2] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu, "Fake News Detection on Social Media: A Data Mining Perspective", Computer Science & Engineering, Arizona State University, Tempe, AZ, USA Charles River Analytics, Cambridge, MA, USA Computer Science & Engineering, Michigan State University, East Lansing, MI, USA
- [3] H. Ahmed, I. Traore, and S. Saad, "Detection of online fake news using n-gram analysis and machine learning techniques," in Proceedings of the International Conference on Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments, pp. 127–138, Springer, Vancouver, Canada, 2017.
- [4] Verónica Pérez-Rosas, Kleinberg Bennett, Alexandra Lefevre, and Rada Mihalcea, —Automatic detection of fake news, Proceedings of the 27th International Conference on Computational Linguistics, pp. 3391–3401, Santa Fe, New Mexico, USA, 2018.
- [5] E. M. Okoro, B. A. Abara, A. O. Umagba, A. A. Ajonye, and Z. S. Isa, —A Hybrid Approach to Fake news detection on social media, vol. 37, no. 2, pp. 454-462, 2018.
- [6] Shlok Gilda, "Evaluating Machine Learning Algorithms for Fake News Detection", IEEE 15th Student Conference on Research and Development (SCORED), 2017.
- [7] Pal, S., Kumar, T. S., & Pal, S. (2019). Applying Machine Learning to Detect Fake News. Indian Journal of Computer Science, 4(1), 7- 12.
- [8] Supanya Aphiwongsophon and Prabhas Chongstitvatana, " Detecting Fake News with Machine Learning Method", CP Journal, 2018.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)