# Fake News Prediction: A Survey

Pinky Saikia Dutta[1], Meghasmita Das[2], Sumedha Biswas[3], Mriganka Bora[4],
Sankar Swami Saikia[5]

[1, 2, 3, 4, 5]Computer Science and Engineering, Girijananda Chowdhury Institute of Management and Technology, Guwahati, Assam, India-781017

**Abstract**— *Fake news generally defined as misleading news often constructed with an aim to create a sense of belief and to mislead people to believe a particular incident. Fake news gets its massive wings through social involvement. We aim to design a system which would probably use concepts like Natural Language Processing (NLP), Data Mining and Machine Learning and prediction classifiers like the Naïve Bayes Classifier and Logistic regression classifier which will predict the truthfulness or fakeness of an article.*

**Keywords**— *Fake news Detection: Logistic Regreesion: Naïve Bayes Classifier: TD-IDF Vectorisation.*

## I. INTRODUCTION

We generally define fake news as something that is verifiably and intentionally false. By false news, we clearly don't consider bias news. News can be biased which depends on the opinion of the person but fake news is intentionally corrupted. The biggest factor behind the success of fake news stories is their high level of social engagement. Social networks connect us with other like-minded people. Our networks of 'friends' on Facebook, or 'followers' on Twitter, generally consist of people who share our values and beliefs. These values may be social, political or economic, and the information we share through these networks helps to define who we are and what we believe in. This identity is then reinforced the more we read similar news stories shared through our social network, confirming our ideas and biases.

As already said fake news is spreading at a very high rate causing social, political and economic problems which are not confined to a small area but is a problem world-wide. These days' fake news is creating different issues from sarcastic articles to a fabricated news and plan government propaganda in some outlets. Fake news and lack of trust in the media are growing problems with huge ramifications in our society. Social media and the internet are suffering from fake accounts, fake posts, and fake news. The intention is often to mislead readers and/or manipulate them into purchasing or believing something that isn't real. So, a system like this would be a contribution in solving the problem to some extent.

As human beings, when we read a sentence or a paragraph, we can interpret the words with the whole document and understand the context. In this project, we teach to a system how to read and understand the differences between real news and the fake news using concepts like Natural Language Processing (NLP), Data Mining and Machine Learning and prediction classifiers like the Naïve Bayes Classifier and Logistic regression classifier which will predict the truthfulness or fakeness of an article.

## II. RELATED WORKS

In this portion, we will relate our work with the existing works.

Regarding fake news detection, the first step is to collect news articles regarding the area we are interested in. for this, we first have to collect a dataset or build it accordingly. The dataset has a set of relevant news articles with true and false labels. This dataset will be the training and testing dataset to train and build the system. All information in the dataset are raw data mostly semi structured and unstructured data. These data are preprocessed for further use.

Natalie Ruchansky, Sungyong Seo and Yan Liu [1] in their journal paper 'CSI: A hybrid Deep model for fake news detection' stated that CSI is a model that combines all three characteristics (i.e. text of an article, user response it receives and the source users promoting) for a more accurate and automated prediction. After incorporating both behavior of the Users and Articles they proposed a model called CSI which is composed of three modules: Capture, Score and Integrate. First two modules based on response, text and sources of an articles using Neural Network to capture the temporal pattern of user on a given article and behavior of users. Based on those two modules third module classify an article as fake or not. This model provides accurate result approximately to 95.3%.

Peter Bourgoje, Julian Moreno Schneider and Georg Rehm [2] in the publication 'From click bait to fake news detection: A approach based on detecting the stance of read lines to article' aimed in detection of the stance of headlines with regard to their corresponding articles bodies and said that the same approach can be applied in fake news, especially clickbait detection scenarios. They took a dataset of classes (unrelated, related, agree, disagree and discuss). First, they checked whether a particular headlines/articles combination is related or unrelated. This is done on "n-gram" matching of the lemmatized input using Core NLP Lemmatizer, 3-class classifier and combined classifier. Best accuracy in related pairs (agree, disagree and discuss) in both classifiers as 79.82 and 89.59.

Manisha Gahirwal, Sanjana Moghe, Tanvi Kulkarni, Devanish Khakhar and Jayesh Bhatia [3] in their publication 'Fake news detection' proposed a system that classifies unreliable news into different categories after computing an F-score using various NLP and Classification techniques to

achieve accuracy. The aim was to accurately determine the authenticity of the contents of a particular news article.

Pérez-Rosas, Verónica & Kleinberg, Bennett & Lefevre, Alexandra & Rada Mihalcea, [4] in their publication 'Automatic Detection of Fake News' focus on the automatic identification of fake content in online news. For this, they introduce two different datasets, one obtained through crowd sourcing and covering six news domains (sports, business, entertainment, politics, technology and education) and another one obtained from the web covering celebrities. They developed classification models using linear sum classifier and five-fold cross-validation, with accuracy, precision, recall and FI measures averaged over the five iterations that rely on the combination of lexical, syntactic and semantic information as well as features representing text readability properties which are comparable to human ability to spot fakes.

E. M. Okoro, B. A. Abara, A. O. Umagba, A. A. Ajonye and Z. S. Isa [5] in their publication 'A Hybrid Approach to Fake news detection on social media' aimed to propose a hybrid model for fake news detection o n social media using a combination of both human based and machine-based approach. Since traditional and machine-based approach have some limitations and cannot single handedly solve the problem like human literacy and cognitive limitations and the inadequacy of machine based approached. To solve all these problems, they proposed a Machine-Human (MH) model for fake news detection in social media. This model combines the

human literacy news detection tool and machine linguistic and network-based approaches. This way two parallel approaches of detection are at work, each helping to provide a balance for the other.

The existing systems and research work reveal that most classification algorithms perform well to detect or predict the fakeness of a news article. Though the logistic regression serves best for the purpose. Our system is based on this information and thus we focus to work with classification algorithms like the logistic regression and a much simpler algorithm like the Naïve Bayes classifier and compare the results of both the classifiers.

### III. EXPERIMENTAL METHOD

As gathered information from the existing systems, the logistic regression algorithm serves best for most the systems. So, we algorithm and thus compare the result of both the algorithms and see which would be ideal for our system.

#### A. Concepts Involved

*Machine Learning*: Machine learning is the field of study that gives computers the capability to learn without being explicitly programmed. It gives the computers the ability to learn which makes it more similar to humans. The basic premise of machine learning is to build algorithms that can receive input data and use statistical analysis to predict an output while updating outputs as new data becomes available.

*Data Mining*: Data mining is the process of analyzing hidden patterns of data according to different perspectives for categorization into useful information which is collected and assembles in common areas such as data warehouse.
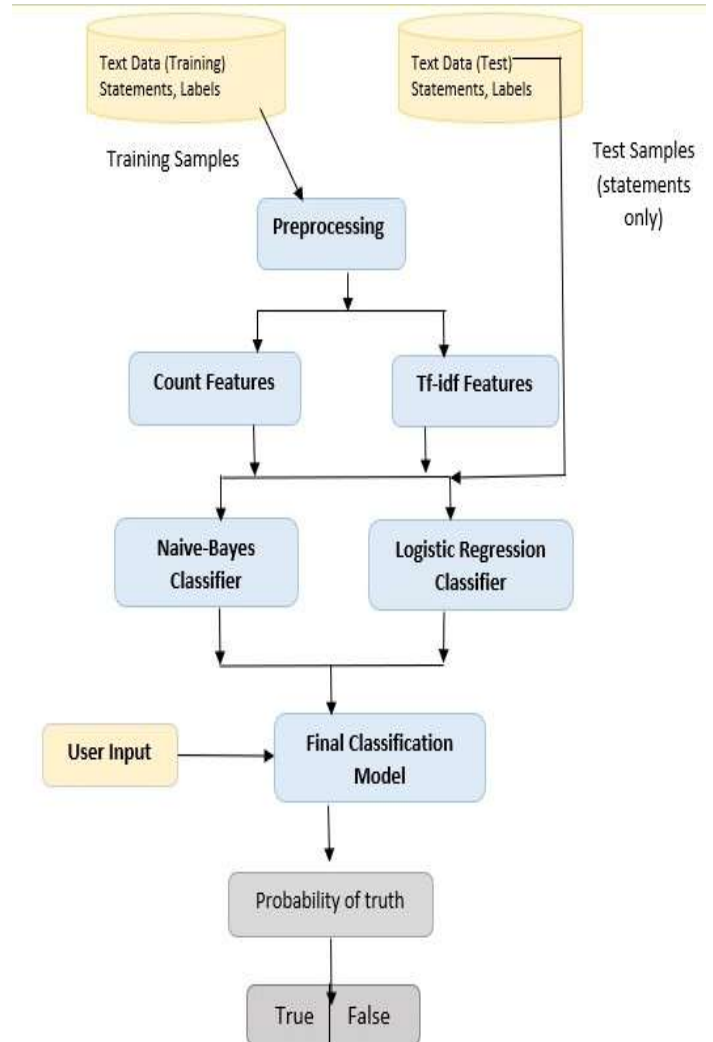


Fig. 1. Flowchart of the proposed system

*Natural language processing*: Natural language processing (NLP) is the ability of computers to understand human speech as it is spoken. NLP helps to analyze, understand, and derive meaning from human language in a smart and useful way.

Data collection is the systematic approach of gathering and measuring information from a variety of sources to get a complete and accurate picture of an area of interest. Data collection enables to evaluate outcomes and make predictions about future probabilities and trends. The name of our dataset is fake_or_real_news and this dataset is available in kraggle.com. This dataset consists of a total of 6337 articles. Then the dataset was manually divided as follows: The training dataset has 80% of the total dataset and the testing dataset has 20% of the total dataset. Now, the training dataset has 5070 articles and testing dataset has 1267 articles.

Data preprocessing is a data mining technique that involves transforming raw data into an understandable format. Real-world data is often incomplete, inconsistent, and/or lacking in certain behaviors or trends, and is likely to contain many errors. Data preprocessing is a proven method of resolving such issues. Data preprocessing prepares raw data for further processing. In preprocessing works like sentence

2

segmentation, tokenization, stop words removal, stemming, lemmatization and removal of stop works is done.

Sentence segmentation is the breaking down of articles into sentences.

Tokenization breaks unstructured data, text, into chunks of information which can be counted as discrete elements. This immediately turns an unstructured string (text document) into a more usable data, which can be further structured, and made more suitable for machine learning.

The words like a, an, the, be etc. these words don't add any extra information in a sentence. Such words can often create noise while modelling. Such words are known as Stop Words.

Stemming helps to create groups of words which have similar meanings and works based on a set of rules, such as remove "ing" if words are ending with "ing".

Lemmatization uses a knowledgebase called WordNet. Because of knowledge, lemmatization can even convert words which are different and can't be solved by stemmers, for example converting "came" to "come".

The process of converting NLP text into numbers is called vectorization in ML. Different ways to convert text into vectors are:

- Counting the number of times each word appears in a document.
- Calculating the frequency that each word appears in a document out of all the words in the document.

Count Vectorizer works on Terms Frequency, i.e. counting the occurrences of tokens and building a sparse matrix of documents x tokens.

TF-IDF stands for term frequency-inverse document frequency. TF-IDF weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus. The mathematical equations for calculating the TF-IFD vectors are:

$$TF(t) = \frac{Number\ of\ times\ term\ t\ appears\ in\ a\ document}{Total\ no\ of\ terms\ in\ the\ document}$$

$$IDF(t) = \log(\frac{Total\ no\ of\ documents}{Number\ of\ documents\ with\ term\ t\ in\ it})$$

Thus,

$$TD\text{-}IDF_{score} = TF*IDF$$

The two prediction classifiers involved in this system are the Naïve Bayes Classifier and the Logistic Regression Classifier.

Naïve Bayes Classifier is a classification model usually helpful with large quantity of inputs. It is a probabilistic classifier based on the Bayes Theorem of Probability with strong independence assumptions. The Naïve Bayes works by using a label with some predefined decisions and thus finds the probability of an unlabeled sample based on the labelled samples provided.

Logistic regression is a predictive modelling algorithm that is used when the Y variable is binary categorical. That is, it can take only two values like 1 or 0. The goal is to determine a mathematical equation that can be used to predict the probability of event 1. Once the equation is established, it can be used to predict the Y when only the X's are known.

## IV. CONCLUSION

In this paper, we will use the Naïve Bayes Classifier and the Logistic Regression Classifier and compare the results of both the classifier. The classifier with the most accurate result will serve as the final model and work with user input. The result will be decided on probabilistic ground and thus the system would predict the truthfulness and fakeness of an article.

## REFERENCES

[1] Natalie Ruchansky, Sungyong Seo, and Yan Liu, "CSI: A hybrid deep model for fake news detection," *CIKM '17 Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pp. 797-806, 2017.
[2] Peter Bourgoje, Julian Moreno Schneider, and Georg Rehm, "From click bait to fake news detection: A approach based on detecting the stance of read lines to article," *Proceedings of the 2017 EMNLP Workshop on Natural Language Processing meets Journalism*, pp. 84–89, 2017.
[3] Manisha Gahirwal, Sanjana Moghe, Tanvi Kulkarni, Devanish Khakhar, and Jayesh Bhatia, "Fake news detection," *International Journal of Advance Research, Ideas and Innovations in Technology*, vol. 4, issue 1, pp. 817-819, 2018.
[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.