

# A Distinctive Approach for Detecting Fake News Using Machine Learning



Md. Rakib Hasan, Ismath Ara Itu

**Abstract:** Nowadays, social media platforms have evolved into an esoteric method for audiences to consume information. News spreading through the social network are also a great source of information as well. For the advancement of Internet access, the information consumption pattern is dramatically changing. As a consequence of those, fake news has become one of the prime concerns because of its potentiality to endanger a society in different perspectives as well as has a political and social impact. Because false news causes so much confusion among people, we will train a model to identify all types of fake news in response to public demand. For this respective research, we collect real data from many reliable and reputed online news portals and fake news from unreliable resources. For converting text to vector format Bag of Words is used. Besides TF-IDF is used for extracting the feature and then CNN is used for classification. With an 83.14% accuracy our model can efficiently detect fake and real news. This work paves a path for an easy automatic fake news detection system which will be very helpful for us to prevent spreading the false information and helps to find the truth.

**Keywords:** Bangla Fake News, CNN, Fake News Detection, Natural Language Processing.

## I. INTRODUCTION

Now-a-days number of online news portal, a great source of information is increasing. News spreading through the social network are also a great source of information as well. For the advancement of Internet access, the information consumption pattern is dramatically changing. As a consequence of those, fake news has become one of the prime concerns because of its potentiality to endanger a society in different perspectives as well as has a political and social impact. Automatic detection of fake news a bit complex because it needs the knowledge of political, social context and some common sense which are lacking in present natural language processing algorithms. In this respective research, a benchmark dataset will be presented for fake news detection which can also be used in fact-checking research and an improved way will be engendered based on different machine learning algorithms and deep neural networks. This work will

pave a path for an easy automatic fake news detection system which will be very helpful for us to prevent the false information and helps to find the truth.

In today's modern world, we are surrounded with different social media platforms, website because of the advancement of technology. From those sites and platforms, we get different news and we consume this news. Most of them people tend to seek out news from this site and platforms than traditional news organization. This traditional organization provide them real news but most of the time we get false news from different social platform and media. Due to this, different machine learning techniques and natural language processing techniques are used to detect fake news. Here machine learning approach is used to perform the classification of real and fake news. And natural language processing is one of the most well-known fields that allows computer to process and manipulate human language. It is used for easy to learn and readable any language to machine. Fake news detection for bangle is more difficult and challenging than other language because of structure of bangle fake news is very confusing. Though some researchers are improving to detect fake news and use different method and algorithms but very few for bangle. There is still a lot of change to improve more.

## II. LITERATURE REVIEW

With the advancement of Information and communication technology more and more online news portal is coming to feed us information. A lot of portals are publishing news without verifying the fact to be the first to publish the news. And off course social media has a big impact for spreading the news. As in social media no one is obliged to share the truth. On the contrary, they sometimes share fake and irrelevant news by choice or just to go with the flow. So, sometimes it's difficult to know which information bears truth and which are not true at all. Many researchers tried to find a way to detect fake news using single or combining multiple machine learning algorithms.

William Yang Wang presents a new, publicly available benchmark dataset named "LIAR" for fake news detection [1]. This database consisting of a huge 12.8k manually labeled short statements is larger than the previously published largest dataset of same type used for fake news detection. They have designed a hybrid CNN to integrate meta data with text that results great in improving a text-only deep learning model.

Manuscript received on January 20, 2022.

Revised Manuscript received on January 25, 2022.

Manuscript published on February 28, 2022.

\* Correspondence Author

**Md. Rakib Hasan\***, Department of Information and Communication Technology, Comilla University, Cumilla, Bangladesh. Email: [rakib@cou.ac.bd](mailto:rakib@cou.ac.bd)

**Ismath Ara Itu**, Department of Management Information Systems, University of Dhaka, Dhaka, Bangladesh. Email: [ismathhasan@gmail.com](mailto:ismathhasan@gmail.com)

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

A fake news detection model is proposed using n-gram analysis and machine learning techniques [2][3]. Two different feature extraction techniques Term Frequency (TF) and Term Frequency-Inverted Document Frequency (TF-IDF) are discussed. Several machine learning techniques are applied to detect fake news and ends with a conclusion of 92% accuracy using Linear Support Vector Machine (LSVM). Some researchers use Count Vector and TF-IDF for extracting feature and then for the detection of fake news they apply machine learning algorithms [3]. Shlok applies TF-IDF of bi-grams and probabilistic context free grammar (CFG) detection on their dataset with different Salgorithms to identify the unreliable sources. They conclude with an accuracy of 77.2% and find a small impact of PCFGs on recall [5].

Lilapati and Rajat work with LIAR dataset from POLITIFACT.COM and concern about the increase of accuracy. They use model ensemble technique to have better accuracy in predicting fake news [6]. Finding gap between how the news is related or unrelated with the real news is the main concern for researchers [7]. Eventually, this approach provides diversity to the field of research with an enviable success rate of 94.21% on test data.

Most of the proposal share their own findings with their accuracy rate for detecting false news. Maximum researchers tried common machine learning algorithms like Support Vector machine, Naïve Bayes, Decision tree, Logistic Regression, Recurrent Neural Network and so on. Some researchers tried their own approaches combining them together but the accuracy is not satisfactory. In these circumstances, an efficient approach is proposed to detect fake news combining deep learning techniques mainly focused on Deep Neural Networks (DNN). This unique approach will provide a new dimension for detecting fake news efficiently with a high accuracy rate as well as lower complexity.

For detecting false news, a supervised system using KNN, Decision Tree, Nave Bayes, Logistic Regression, and SVM works well. However, owing to the large dataset required, development will be time consuming. The amount of the training data determines its performance level. Furthermore, the neural network has done a good job at this, but it is still difficult to improve. It is now a word for a system.

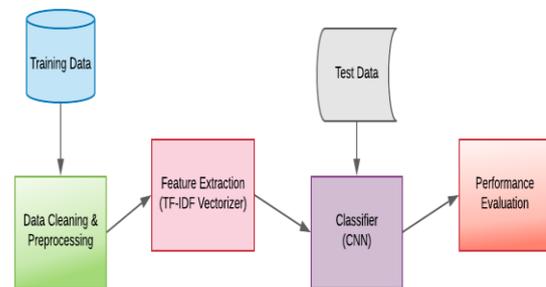
### III. METHODOLOGY

The aim of this research is to determine the performance of a system in the identification of false news. The following actions were followed to complete the thesis effectively.

- To build dataset collect data from different reputed and non-reputed news portal by scraping using different scraping tools including (ParseHub, Scrapy etc.)
- Collect the data of different categories.
- Preprocess the data.
- Split the data for train and testing.
- For data read and write panda framework is used.
- Data.shape and data.columns to show dataset effectively.
- Check missing data. For this we have used isnull() function.
- Observe the split data as real data and fake data of dataset.

- Count all the categories.
- Count the data according to categories.
- Fill the missing data with true value.
- Identify, reduce noise and remove all type of noise.
- Target variable encoding to create data shape and partition data.
- Tokenize the data.
- Using Bag of Word, you may transform text documents into numerical characteristics. The simplest is the Count Vectorizer, which counts the number of times a token appears in the text and uses that value as its weight. Finally, to improve accuracy, an unsupervised Bangla POS tagger based on suffix analysis is suggested.
- By using tf-idf convert numeric data into matrix format and
- Finally, construct training and testing sets then model train using CNN.

The mentioned steps are the detailed description of this thesis. But in a nutshell, they can be categorized in some main steps.



**Fig. 1. System flow Diagram of the Proposed System**

In Figure 1 the main steps are depicted. Firstly, we need to collect data. Then preprocess the data is a must for better outcome. After preprocessing, the data is ready for feature extraction. Using those feature the model is implemented and then testing the model with test dataset. Finally, it comes to performance analysis. These are main steps in short.

#### A. Data collection Procedure

Data collection mainly includes data acquisition, data labeling, using existing data. Different types of tools and software are used for different tasks like data discovery, data augmentation, data generation, data labeling, improving the existing data.

For performing those tasks, we can use different scraping tools (Scrapy, BeautifulSoup, Scraper, ParseHub, Octoparse, Mozenda etc.) which will be used for collecting data from web, OpenCV and scikit-image as data augmentation tools. Lionbridge AI will be used for handling image, text, video data and LabelBox will be used for image annotation, text classification.

Our data is split into six columns (Title, Date, Category, Headlines, Description, Class) and classes are classified into two categories Real and Fake.

Our entire dataset is about 2082, with 80% of the data being utilized for training and 20% for testing. Data from reputable Bangladeshi online news portals such as Daily Prothom alo, Bangladesh Pratidin, Ittefaq, Daily KalerKantho, Daily NayaDiganta, bdnews24.com, and others was collected.

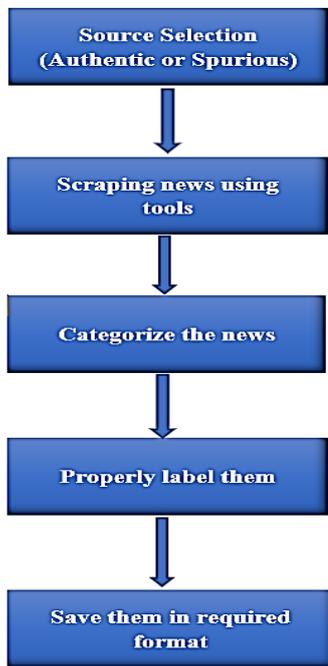


Fig. 2. Data Collection Procedure

Fake class data was gathered from unreliable internet news sources such as Dhaka Channel, baler kontho, moti kontho, Khbor24.com and others. After that, we use Google scraping to verify whether or not the data is false. For web scraping, we must search by news title, create a result for that news, and determine how much of that news is shown in a fake news portal vs a genuine news portal. But this process is very time consuming and hectic. So, we have used some paid scraping tools like ParseHub, OctoParse to web scraping. In ParseHub, at first, we need to download the application from their website. Then need to open a new project. Provide the website url you want to collect data. Then from the user interface one can select the title and description of news as well as any other information. After the “Get Data” and “Run” one can have the option of downloading the extracted news in required format. We divided the news into two categories based on the aforementioned description: false news and genuine news. Figure 2. depicts several major topics as well as a flowchart depicting how we gather and identify data. In Figure 2, we also indicate which sites provide genuine data and which sites provide false data. As a result, this flowchart is split into two parts for news sites, one reputable and the other not. We also observed a news category while gathering data, which indicates which news was produced for a joke and which news was published for an entitlement.

We sometimes get news with an appealing headline, but when we read the item, we notice something else. To effectively identify news, one must analyze all news, which is referred to as reading beyond the headlines. Here are some important aspects we have considered to collect data properly:

- To collect valid data, we consider the reputed source for authentic data and collect fake data from non-reputed source.
- We must consider the full news. Only headline cannot comment actual scenario.
- To identify data that’s the data is real or fake must be check author because some non-reputed author publishes fake news and reputed author always publish real news so it’s so important to detect data and collect valeted data.
- To collect news, we focused on news date which claims more authenticity.
- At the time of collecting data, we checked this data is type because it’s helpful to arrange data into different categories.
- For data validation, we recheck the news on Google by hold news title by use news title we search Google and find some result if this news is real then show some reputed news site and if this news is fake then show non-reputed site or only show one site because actually fake news published non-reputed site and if is news is fake then it published individual site.

Figure 3 depicts a sample of authentic data. We have arranged 2082 dataset; however, we will only present some datasets as examples. We store news headlines, categories, contents, domain, and classes in our dataset (classes indicate which news is genuine and which is false; we have two classes, fake (for fake news) and real (for true news) (for real news).

Label	1 (Authentic)
Domain	Somoynews.tv
Published Date and Time	9/19/2018 19:37
Category	Sports
Headline	জাতীয় দলের একটু কাছে আশরাফুল
Contents	খুলনায় কাল শুরু হচ্ছে একটি চারদিনের ম্যাচ। এই ম্যাচে থাকছেন আশরাফুল। নিষেধাজ্ঞা কেটে যাওয়ার পর তিনি প্রথমবারের মতো সুযোগ পেয়েছেন বিসিবি'র কোনো দলে খুলনায় শেখ আবু নাসের স্টেডিয়ামে কাল শুরু হচ্ছে একটি চারদিনের ম্যাচ। এই ম্যাচে খেলছেন মোহাম্মদ আশরাফুল। ম্যাচটা আশরাফুলের কাছে বিশেষ তাৎপর্যময়।

Fig. 3. A Sample Authentic Data

Here we show two different datasets for real and fake news. Real news dataset is shown and leveled as 1. The real news is collected from the reputed news sources on the other hand the fake news is collected from unreliable and non-reputed sources and leveled as 0.

As we mentioned, we have worked with 2082 data. In them 1500 data are collected from reputed sources and considered as real news. Rest of the data are considered non reliable and fake. Figure 4 shows the list of fake news.





Table II depicts the tokenized version of a raw sentences or string.

**Table-II: Tokenized Version of A Sample**

<b>Raw sentence</b>	রোহিঙ্গা নির্যাতনে আন্তর্জাতিক অপরাধ আদালতের প্রাথমিক তদন্ত শুরু
<b>Tokenized form</b>	“রোহিঙ্গা” “নির্যাতনে” “আন্তর্জাতিক” “অপরাধ” “আদালতের” “প্রাথমিক” “তদন্ত” “শুরু”

The bag of words model is one of a number of methods for extracting characteristics from text that come from the area of computer science known as Natural Language Processing, or NLP [10]. It does this by counting the number of times words appear in a document. The majority of the time, a bag of words is utilized for numeric format. The same word appears many times in a single line. Numeric format allows us to count their frequency.

In Table III, “রোহিঙ্গা” occurs 3 times from this text script, then “আন্তর্জাতিক” occur 3 times, “মিয়ানমার” occur 2 time, “বাংলাদেশে” occurs 2 times. Finally using of bag of word process, we get frequency of word. And we can simulate this text script and represent numeric format. In this process we have to import pickle library. This library used for vectorization. CountVectorizer() method is used for frequency define with some required arguments (max\_features=1500, min\_df=4, max\_df=0.8) . Here df means that document frequency, max-df remove too much frequency. Here max df =.8 means that more than 80 percent frequency it will remove. And min frequency remove number of low frrequency. Here min\_df=4 that means less than 4 frequency it will remove.

**Table-III: Frequency Count of A Sample News**

রোহিঙ্গাদের বিরুদ্ধে মানবতাবিরোধী অপরাধে অভিযুক্ত মিয়ানমার সরকারের বিচারে প্রাথমিক তদন্ত শুরু করেছে আন্তর্জাতিক অপরাধ আদালত। এর মধ্য দিয়ে রাখাইনে সেনা অভিযান নিয়ে, একটি পূর্ণ তদন্তের পথ সুগম হলো। রোহিঙ্গাদের দেশত্যাগে বাধ্য করা, গণহত্যা, ধর্ষণের পাশাপাশি মৌলিক অধিকার হনন ও নুটপাটের অভিযোগের তদন্তে নেমেছে আন্তর্জাতিক অপরাধ আদালত, আইসিসি.....	Token	Frequency
	“রোহিঙ্গা”	3 Times
	“আন্তর্জাতিক”	3 Times
	“মিয়ানমার”	2 Times
	“বাংলাদেশে”	2 Times
	.	.
	.	.
	.	.

For translating text to numbers, the bag of words method works well. It does, however, have one flaw.

It gives a word a score depending on how many times it appears in a document. It doesn't account for the fact that the term may appear often in other texts as well. TFIDF addresses this problem by dividing a word's term frequency by its inverse document frequency. The letters TF and IDF stand for "Term Frequency" and "Inverse Document Frequency," respectively. [9].

The term frequency is calculated as:

$$\text{Termfrequency} = \frac{x}{z} \dots\dots\dots (I)$$

X =Number of Occurrences of a word

Z =Total words in the document

And the Inverse Document Frequency is calculated as:

$$\text{IDF} = \left[ \frac{V}{N} \right] \dots\dots\dots (II)$$

V =Total number of documents

N =Number of documents containing the word

Every word is calculated in a particular way. Using this method, we get tf and idf value then it is converted into array. To prepare our dataset for the model, we first need to extract the feature using this tf-idf method. We have applied some methods to accurately calculate the TD-IDF and then fit the dataset for the model.

**E. Classifier**

Before Classification we need to perform the tasks including data processing, data cleaning, feature extraction as mentioned earlier. To feed the data into classification algorithms train and test dataset was gained by splitting the main dataset in two different sets. The splitting ratio was 80% for the training dataset and 20% for the test dataset. Text classification or sentence categorization into a set of predefined categories using n-grams, i.e., words or characters, sequences of words or sequence of characters is considered one of the most common tasks in natural language processing. ConvNet are employed to perform the task of classification and make the classification task easier. For this, a series of words w<sub>1</sub>, w<sub>2</sub>, w<sub>3</sub>, ..... , w<sub>n</sub> is needed to feed to the convolutional network. We get a vector from bag of words that can also be feed to CNN.

Each of the vector is connected with a d-dimensional embedding vector. If a sliding-window of size k moves over the sentence we get a one-dimensional convolution of width k. Every window of the sequence is subjected to the same one-dimensional convolution filter, also known as a kernel. Particularly, a dot-product of the concatenation of the embedding vectors in a given window and a weight vector, followed by a non-linear activation function.

The data is pooled into a single one-dimensional vector in ConceptNet, which is a network of convolution windows.

The most essential elements of the sentence/document should be included in this vector. The notion that ConvNet is merely a feature extractor, leading to a fully connected layer for prediction.



**Fig. 5. Real News Word Cloud**





9. Hadeer Ahmed, Issa Traore, SherifSaad, "Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques, Conference paper, First Online: 11 October 2017.
10. "Bag of words model", *Wikipedia*. Accessed on: May 05, 2021. [Online]. Available: [https://en.wikipedia.org/wiki/Bag-of-words\\_model](https://en.wikipedia.org/wiki/Bag-of-words_model)
11. "Tokenization in NLP", *Analytics Vidhya*. Accessed on: April 20, 2021. [Online]. Available: [www.https://www.analyticsvidhya.com/blog/2020/05/what-is-tokenization-nlp/](http://www.https://www.analyticsvidhya.com/blog/2020/05/what-is-tokenization-nlp/).
12. "What is Anaconda and how does it relate to Python", *VENTURE LESSONS*. Accessed on: Jan 15, 2021. [Online]. Available: <https://www.venturelessons.com/what-is-anaconda/>

### AUTHORS PROFILE



**Md. Rakib Hasan**, has completed his B.Sc. (Honors) and M.Sc. in Computer Science and Engineering from Jahangirnagar University, Dhaka, Bangladesh in 2016 and 2018 respectively. He is now working as a Lecturer at Department of Information and Communication Technology, Comilla University, Cumilla. Previously he worked as a Lecturer at

Department of Computer Science and Engineering, Daffodil International University, Dhaka from May 2016 to February 2019. His teaching experience includes different graduate (M.Sc.) and under graduate courses. He is currently working on Artificial Intelligence and Machine Learning. His research interest includes Natural Language Processing, Machine Learning, Cyber Security and Computer Vision. He is actively engaged in educational activities.



**Ismath Ara Itu**, has completed her B.B.A and M.B.A. in Management Information Systems from university of Dhaka, Dhaka, Bangladesh in 2015 and 2016 respectively. She is currently working on Data Science, Secure Information Systems and Machine Learning. Her research interest includes Machine Learning, Data Science and Information System Security.