**Topic:** **Fake News Detection Using Machine Learning**

**Submitted by**

**Hasibul Islam**
**ID: 191-35-2685**
**Batch: 28**
**Department of Software Engineering**
**Daffodil International University**

**Supervised by**

**Mr. Md. Rittique Alam**
**Lecturer**
**Department of Software Engineering**
**Daffodil International University**

The report for the Bachelor of Science in Software Engineering have been met by the submission of this thesis paper.

# APPROVAL

This thesis titled on "FAKE NEWS DETECTION USING MACHINE LEARNING", submitted by **Hasibul Islam (ID: 191-35-2685)** to the Department of Software Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of Bachelor of Science in Software Engineering and approval as to its style and contents.

**BOARD OF EXAMINERS**

------------------------------------------------  **Chairman**

**Dr. Imran Mahmud**
**Head and Associate Professor**
Department of Software Engineering
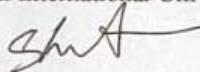Faculty of Science and Information Technology
Daffodil International University

------------------------------------------------  **Internal Examiner 1**

**Md. Khaled sohel**
**Assistant Professor**
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University

------------------------------------------------  **Internal Examiner 2**

**Md. Shohel Arman**
**Assistant Professor**
Department of Software Engineering
Faculty of Science and Information Technology
Daffodil International University

------------------------------------------------  **External Examiner**
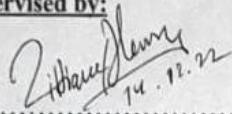
**Rimaz Khan**
**Managing Director**
Tecognize Solution Limited
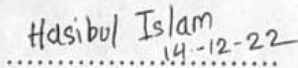
# Thesis Declaration

I announce hereby that I am rendering this study document under **Md. Rittique Alam, Lecturer,** Department of Software Engineering, Daffodil International University. I therefore state that this work or any portion of it was not proposed here therefore for Bachelor's degree or any graduation.

**Supervised by:**

..............................

Md. Rittique Alam
Lecturer
Department of Software Engineering
Daffodil International University

**Submitted by:**

..............................

Hasibul Islam
ID:191-35-2685
Department of Software Engineering
Daffodil International University

# ACKNOWLEDGEMENT

*I express my gratitude to Almighty Allah for giving me the ability to complete the thesis.*

*I am especially thankful to my supervisor Md. Rittique Alam (Lecturer)Faculty of Science and Information Technology, Department of Software Engineering, Daffodil International University . He has helped me in different ways at different times. Whenever I was faced with obstacles, he provided maximum assistance to complete the thesis.*

## Hasibul Islam

# ABSTRACT

Today, fake news detection is a trendy topic that has drawn a lot of attention from academics around the globe. Any content that is untrue and created with the intention of leading readers to believe a lie is typically considered fake news.

In this paper, a framework is proposed that should be used to begin the project. It calls for the application of text-processing, cleaning, and feature-extraction techniques to reorganize the information, which should then be "obeyed" into each classification model during training and parameter tuning to produce the most accurate and optimized predictions for identifying fake news.

This study examines three example datasets to better understand the background for identifying fake news. It also makes an effort to determine linguistic differences between false and real news items using a variety of visualization techniques. This text's goal is to provide a detailed analysis of the results of several popular machine learning classifiers, including the Support Vector Machine, the Naive Bayes Method, the Decision Tree Classifier, the Random Forest, and the Logistic Regression, as well as the development of the Ensemble Method (Bagging & Boosting), which uses classifiers like the XGBClassifier and the Bagging Classifier to combine various amounts of classification models for identification.

Keywords:- Detection Fake News, Scraping, Social Media, Text classification, Comparison of Algorithms, Machine learning, Natural language process.

# Table of content

# Table of Figure

# List of Tables

**CHAPTER 1**

INTRODUCTION

Today, fake news detection is a trendy topic that has drawn a lot of attention from academics around the globe. Any content that is untrue and created with the intention of leading readers to believe a lie is typically considered fake news.

To get the most out of the optimized and best predictions for the fake news, the proposed framework in this paper should be "observed" in each of the training and tuning parameters classification models for each model. This is done by applying textprocessing, cleaning, and features extraction strategies at the beginning of the project, which are meant to rearrange the content.

This study attempts to break down the linguistic differences between fake and real articles in order to include visualization of findings using a variety of visualization tools. It uses three use case datasets for the proposed system to understand the percentage of data that is responsible for identifying fake news. This text's objective is a thorough evaluation of the results of several popular machine learning classifiers, including the Support Vector Machine, the Naive Bayes Process, the Decision Tree Classifier, the Random Forest, and the Logistic Regression, as well as the development of the Ensemble Method (Bagging & Boosting), such as the XGBClassifier, the Bagging Classifier of various prediction combinations of classification models, to determine which provides the best optimal results.

## 1.1 Problem Definition

Identifying fake news on social media presents a number of difficult research issues. First, fake news is purposefully crafted to confuse readers, making it difficult to distinguish from merely supported information.news material As a result, linguistic traits alone are inappropriate for identifying fake news.Second, to spice up detection, specific auxiliary information types of content and user social engagements must even be supplied. However, the exclusion of this auxiliary information leads to a unique and significant problem with the sensitivity of the data. Despite the fact that information from various modalities can offer hints for fake detection, it can be difficult to draw out distinctive traits from each modality and successfully combine them. Formally, our issue is as follows:

We may want to compute and comprehend whether the news articles are fake news or not in the pretend news detecting problem. In our dataset set, we'll represent the label as Y =, where 1 represents "Fake News" and 0 represents "Real News." Predicting whether or not article A might be a fake news article is the goal of the pretend news detection framework

$$f(a) = \begin{cases} 1, & \textit{it occure if a take the role real news} \\ 0, & \textit{Otherwish this is fake news} \end{cases}$$

# 1.2 Motivation

Through my research on fake news and the accompanying thesis papers, I have come to the conclusion that fake news may also have a negative impact on society, how we view the media, and how we perceive facts and opinions in general. If we want to avoid reality-induced dizziness and safeguard our society, especially the less educated members of it, it is imperative that bogus news be identified and classified as either fake or true.

In this context this text makes the next contributions:
 • Prevalence of faux news on social media
 • Emerging research area in tongue Processing (NLP).
 • Will Evaluation environment, competitors, datasets, performance measures to ready to i will be able to be able to find the only accuracy using ML models
 .• Basic countermeasures are inflexible and inefficient.
 • Data miss understands to processing for best evaluation and prediction
 • To avoid taken prevalence to spread the fake news.

## 1.3 Thesis Contribution

The main contribution of this research could also be summarized as follows:-

I've worked hard to produce the following contribution in this thesis paper. In the end, this research paper's goal is to use a scientific report from a dataset of fake news to help academics better understand which model and technique are the most effective. This report will also aid in the implementation of the tools or models in the future. This paper used a variety of strategies to improve both model performance and accuracy. Our research's primary goal is to strengthen the model so that it can find the most accurate estimate result from the training dataset. On a single, divided standard dataset, we also compare the results of seven to eight machine learning (ML) algorithms and ensemble techniques, such as Support Vector Machine, Nave Bayes Method, Decision Tree Classifier, Random Forest, Logistic Regression, and XGBClassifier.

## 1.4 Thesis Organization

The remaining sections of this thesis paper are as follows:

It was chosen since the inquiry on false news identification uses the shown in Section 2. With experiments where we compared different confusion matrices of the method to see which one could best suit our model using many Python libraries, we may examine performance matrix, accuracy, precision, recall, and F1-score, which are briefly presented in Section 3. On the other hand, the datasets, the features, the experiments likewise because the results of our experiment on three-part of Case Study are summarized in Section-4.Finally, it has been proposed that Gradient Boosting Algorithm (XGB) may efficiently recognize fake news with extremely high accuracy and F1 score using the right set of characteristics collected from the texts and headlines, as detailed in Section 5.

## 1.5 Scope

The scope of this thesis are exiting helpful for:-
 • Detector Identify the True/Fake news or Article among Fifty Thousand data to help avoid the rumor.
• Investigate could also be an extended process to identify which is true or false news but this method will help to identify fake news within a second.
 • Prediction supported "text or article" using machine learning will help people, politicians, and industrial level, etc.
 • Providing the only prediction will help people in touch in mind of the rumor.
 • Applying the Ensemble classification achieved the foremost effective accuracy which clearly said this work will provide trusting information which is most are helpful for the politicians

# CHAPTER 2
## LITERATURE REVIEW

### 2.1 Introduction

I'll expand on some similar works that are interesting to look into in this section. Recently, a number of approaches to the ubiquitous problem of fake news and aim to as accurately differentiate fake information

### 2.2 Comparative Performance of Machine Learning Algorithms for Fake News Detection

Seven machine learning (ML) algorithms are compared in this paper by Arvinder et al.using three common datasets, a completely original set of features, and statistically valid the outcomes using accuracy and F1 ratings. There has been a claim made in the conclusion section that Utilizing a proper collection of features taken from the texts and also the headlines, Gradient Boosting Algorithm (XGB) can successfully recognize bogus news with very high accuracy and F1 score.

**Many hand-crafted features have been used to feed this network, such as**:-
 • N-grams Count Feature: These features are used for counting occurrences of n-grams within the title and body of the news, and various ratios of the unique n-gram and total word count, and this n-gram could also be unigram, bigram, or trigram.
 • Word Embedding: Forgetting the vector space representation of the words, they used the Word Embedding technique. Word Embedding replaces each word with a real-valued vector.
• Sentiment Polarity Score: for creating the good choice as a feature used SP score for locating the sentiment of the news.
 • Psychological Features: Create using Linguistic Investigation and Word Count, which can be a comprehensive dictionary developed by text mining tools.
 • Feature Matrix: Features and their counts are summarized in figure 1

| Features | | Number of feature vectors |
|---|---|---|
| Sentiment | Headline | 4 |
| | Content | 4 |
| Readability | | 12 |
| Count | | 41 |
| Cosine Similarity of Normalized *tf-idf* vectors between headline and content | | 1 |
| Word Embedding | Headline | 50 |
| | Content | 50 |
| | Cosine similarity between Headline and Content | 1 |
| Total number of features | | 163 |

**Their outcome are shown in figure 2**

| Feature excluded | | Dataset 1 | Dataset 2 | Dataset 3 |
|---|---|---|---|---|
| tf-idf | F1-score | 0.9 | 0.94 | 0.89 |
| | Loss | 0.02 | 0 | 0 |
| Count | F1-score | 0.89 | 0.92 | 0.88 |
| | Loss | 0.03 | 0.02 | 0.01 |
| Word embedding | F1-score | 0.8 | 0.9 | 0.85 |
| | Loss | 0.12 | 0.04 | 0.04 |
| Sentiment | F1-score | 0.9 | 0.94 | 0.89 |
| | Loss | 0.02 | 0 | 0 |
| Readability | F1-score | 0.9 | 0.93 | 0.88 |
| | Loss | 0.02 | 0.01 | 0.01 |

## 2.3 Which machine learning paradigm for fake news detection?

In order to protect people, especially those who are less educated, from the purported reality vertigo, Dimitrios et al. [6] developed their solution. To address this problem, various AI strategies

are suggested.An detailed presenting evaluation of eight AI computations for identifying and characterizing fake news is provided in this article.

The computations for the test are listed below and include L1 Regularized Logistic Regression, C-Support Vector Classification, Gaussian Naive Bayes, Multinomial Naive Bayes, Decision Trees, Random Forests, Multi-Layer Perceptrons, and Convolutional Neural Networks. Sci-pack Learn's version of the first seven calculations is the one to use, however they came up with their own formula for the eighth calculation.

## The result is shown in figure 3

| Model | 50D | 100D | 300D | Glove Vectors 50D | 100D | 300D | TF/IDF |
|---|---|---|---|---|---|---|---|
| LR | 0.69-0.01 | 0.97-0.01 | 0.58-0.01 | 5.75-0.01 | 7.36-0.0 | 3.13-0.01 | 0.04-0.0 |
| MLP | 8.37-0.0 | 7.4-0.0 | 11.45-0.0 | 8.12-0.0 | 6.45-0.0 | 10.74-0.0 | 8.46-0.0 |
| DT | 1.1-0.0 | 02.01-0.0 | 6.39-0.0 | 1.1-0.0 | 1.76-0.0 | 5.44-0.0 | 0.58-0.0 |
| RF | 1.02-0.01 | 1.39-0.01 | 2.31-0.01 | 0.96-0.01 | 1.33-0.01 | 2.26-0.01 | 0.84-0.01 |
| GNB | 0.01-0.0 | 0.01-0.0 | 0.03-0.01 | 0.01-0.0 | 0.01-0.0 | 0.03-0.01 | 0.05-0.01 |
| MNB | 0.01-0.0 | 0.01-0.0 | 0.03-0.0 | 0.01-0.0 | 0.01-0.0 | 0.02-0.0 | 0.00-0.00 |
| SVM | 14.44-01.08 | 19.08-1.68 | 54.86-4.81 | 13.04-01.09 | 19.09-1.72 | 53.44-4.78 | 10.39-0.91 |
| CNN | 9.88-0.24 | 12.28-0.27 | 16.99-0.27 | 12.11-0.29 | 14.72-0.28 | 17.15-0.29 | |

**Figure 3: Training/classification times (in seconds) for Dataset1.**

## 2.4 Hybrid Machine-Crowd Approach for Fake News Detection

About automated deceit detection, there are several important articles. In [7], writers used the hybrid machine-crowd approach as a good way to deal with the problem of false information in general. The author specializes in identifying satire or parody and fake content using the public Fake vs Satire dataset; this method provides improved accuracy at the right price and latenc .This method combines the efficiency of system learning algorithms with the knowledge of crowds, using crowdsourcing in the situations where system learning algorithms fall short of performing with high accuracy. As a result, this method gives improved accuracy at an appropriate price and latency. This approach is simple enough to be easily applied to various datasets and experimental setups for concerns with false information detection.

The following is a summary of their contributions:

• They find vital clean to compute functions that enhance the baseline accuracy through 2.54% at the computerized detection the usage of system gaining knowledge of (81.64% as compared to 79.1%).

• They have designed a crowdsourcing undertaking that leverages the fact-checking capabilities of on-line crowd employees to choose information content material, attaining an accuracy of as much as 84%.

• They layout a hybrid faux information detection machine as a trade-off among accuracy, latency, and price. The hybrid method will increase the overall accuracy through as much as 87%

**The proposed Hybrid Machine-Crowd Approach as follows:-**



**Figure 4: The hybrid machine-human concept designed for fake-news detection**

**The result they have found to apply this approach as follow the figure 5**

|  | fold- 1 | fold-2 | 10 folds |
|---|---|---|---|
| **Top k models**<br>**Prob.threshold**<br>**voting** | 3<br>0.72<br>2 | 3<br>0.72<br>2 | 3<br>0.72<br>2 |
| **accuracy** | 96% | 84% | 86.79 |
| **% of tasks to crowdsource** | 34 | 41.5 | 42 |
| **Estimated cost** | $2.7 | $3.4 | $34 |
| **Estimated latency** | 1h | 1.3h | 28h |

**Figure 5: Results from the hybrid approach**

## 2.5 FNDNet: A deep convolutional neural network for fake news detection

To identify misleading information, Rohit et al. [8] proposed the deep convolutional neural network (DNFNet) version. Their variant (FNDNet) has been developed to automatically evaluate

discriminating functions for the category of false information through a few deep neural network layers that are hidden from view. It is not programmed to rely on the created functions instead. For each layer, it creates a deep neural network (CNN) to eliminate a number of functions. They examine the overall effectiveness of the suggested system using several baseline models. By examining the specifics, benchmark data sets were employed to educate the study of the edition, and the suggested text also attained 98.36 percent of the kingdom of the effects of art. Several broad criteria are used to validate the results, including

**The result they have proposed using machine learning and deep learning-based models shown in figure 6:-**

| Word Embedding Model | classification model | Accuracy(%) |
|---|---|---|
| Tf-Idf on unigrams and bigrams | Neural Network | 94.30 |
| BoW without unigrams and bigrams | Neural Network | 89.22 |
| word2Vec | Neural Network | 75.68 |
| GloVe | Mutinomial Naive Bayes | 89.95 |
| GloVe | Decision tree | 73.66 |
| GloVe | Random Forest | 71.33 |
| GloVe | KNN | 53.70 |
| GloVe | CNN | 91.55 |
| GloVe | LSTM | 97.30 |
| GloVe | our propose model | 98.37 |

**Figure 6: Results of Classification using Machine Learning and Deep Learning-based models.**

## 2.6 Detecting Fake News using Machine Learning and Deep Learning Algorithms

Tanvir et al.'s[12] recommendations for machine learning and deep learning techniques to identify incorrect knowledge are numerous. This uses a variety of categorization models, such as Help Vector Machine, Naive Bayes, Logo Regression, Long Immediate Memory, and Recurrent Neural Network. Additionally, combining these categorization methods enhances estimation even more. These models can be used to assemble a crossvalidation from a training dataset using k-fold (k=2) using the Sci-pack, predicting with data collection.They looked at a computerized model for Twitter post authentication that offers broad solutions and examples of news that has been misidentified for the purpose of information aggregation. A deep learning algorithm is suggested to identify incorrect information once a method has been constructed from the monitored models.

They choose a number of positions for success observation using distinctive, regulated, and profound learning approaches. There are four usable vectors that are eliminated.

• **Count-Vectorization**

• **Character level-vectors**

 • **Word Level Vectors**

• **NGram vectors**

The Vector Machine support was carried out to track the previous findings to some degree further. However, no changes have been found at this point as 74 percent of the four function vectors cited in their analysis are of research accuracy. The above results are seen in the figure.



**Figure 7: Accuracy Comparison**

Any of the vector features of the Naive Bayes model listed above must be evaluated on them. For 73% of the vector number functions, 75% of Word Level TFIDF, N-grams, and character vector, more precision is provide

Finally, a logistic regression model was used. This has marginally boosted production compared to before, with vector count and word ratios predicting 74% and 76%, respectively. The results of the logistic regression ranged from 75% to 76%.

**The final result they have shown as follow in figure 8**

|  | Precision | Recall | F1-Score |
|---|---|---|---|
| **Naive Bayes** | **.89** | **.99** | **.94** |
| **Logistic Regression** | **0.89** | **0.75** | **0.81** |
| **SVM** | **0.89** | **1.0** | **0.94** |

**Figure 8: Score of Recall, Precision and F1-Score for Naïve Bayes, Logistic Regression, SVM**

## 2.7 Conclusion

The majority of such works, as we have seen in the past, focus on enhancing prediction accuracy through the addition of new attributes. The fact is that these features do not always exist; for instance, some periodicals do not have images. Utilizing information from social media is also extremely dangerous because it is simple to create a new account on these platforms and fool the identification system. Because of this, I only chose to consider the article's content in order to determine whether misleading

# CHAPTER 3
# METHODOLOGY

The aforementioned subject teaches us that there are two categories for news, whether it is accurate or incorrect. Before selecting our model and evaluating the effect on a blood type test, it is important to understand the significance of the issue. AI is lots of calculations, but some of them are useful in determining "Reality or Fiction" and some are on a daily basis. Our main focus was feature engineering so that the accuracy of recognizing the news might be a great lot of talent if we could tweak the feature or add a different feature. From the perspective of a mental discovery of fake news, we learn that the word count in an exceedingly documentary expression

## 3.1 Collect Dataset

We begin by selecting a Data-set that can be utilized to distinguish between Real News and Fake News. As was already mentioned, the main focus of this was a newspaper that dealt with Fake News strands about our culture.our nation, both politically and internationally. This offers details about (mainly utilized) political, social.In the beginning, fake news datasets were gathered from Kaggle[1], where this dataset relates to fake news spreads during the 2016 U.S. Presidential Election. Additionally, IEEE Dataport and some of the True news have been gathered using web scraping from reliable and relevant sources like CNN, BBC, and The New York Times, among others.

Finally, data that have already been gathered from many sources and combined with fake and real data have been created. This dataset was utilized in this study to test the model we suggested. Finally, fake and true data has been concatenated which are already collected from different sources. This paper has been used this dataset for our proposed model. There is a total of 44909 (https://doi.org/10.6084/m9.figshare.13325198.v1) data in the dataset. The total dataset has been divided into three case studies of data. 44,909 document has been assign as CS-1, 33,681 documents have been assign as CS-2 and 22,454 documents has been assign as CS-3 among total dataset.

For selected 100% document the calculation is:


**CS-1 =(100/100)\*Total Number of Document**


For selected 75% document the calculation is:

**CS-2 =(75/100)\*Total Number of Document**


For selected 50% document the calculation is:

**CS-3 = (50/100) \*Total Number of Document**


This data collection includes F and T mark groups, where F stands for fake news and T for true news, indicating what is tagged as Actual or Fact. The dataset also includes a single column for statements that include all linkages and information. There are four attributes: Title, Text, Topic, and Date. And the information was labeled as Real News and False News.



**Figure 9: Divided to the subject level function in various groups**


As Illustration above Fig 9, that's represents the subject feature. This subject feature has four classes. Where in 'news' group has 29.1% data, 'politics' have 40.3% data, 'world news'' have 29.3% data and 'US_news' have 1.29% data among 100% of our dataset.

## 3.2 Data Preprocessing

After a diversification of knowledge, the arrangement and reduction of information in a form that is effectively decoded and understood by computers becoPreparing knowledge for a machine learning model through data preprocessing is another possible way. This is both the first and most important stage in creating a machine learning model. We don't always come across clear, well-formatted data while creating a machine learning project. Additionally, every time you handle data, you must clean it up and arrange it perfectly. But for this, we're employing pre-processing tasks.mes vitally necessary in the field of machine learning.

Preprocessing has attested when data are inconsistent or have not transformed all data to lowercase as though the information does not abuse the machine at the beginning of awareness. After that, we disabled a feature that wasn't necessary for identifying phony news. The best project

Then, Daffodil International University undertook to eliminate the pointing markings that were used to print and write different sentences and clauses and to clarify the concept of sentences. As stop words, commas, periods, and question marks are examples of punctuation that are irrelevant because they cause noise in the dataset. We discovered the overall number of process study datasets after preprocessing the specifics. The pre-processing information technique described below.



Figure 10: Framework of FakeSpy for training and testing algorithms to classification of news articles

## 3.3 Features Extraction

When dealing with a large amount of data, features extraction techniques should be used because the majority of them are repetitive and irrelevant, making calculations difficult and producing false results. Feature extraction is also a general term for putting together a set of factors to exacerbate these problems while still portraying information through advertising. Many experts in artificial

intelligence (AI) and machine learning assume and concur that correctly enhanced component extraction results from successful model construction.

# 3.3.1 Term Frequency-Inverse Document Frequency (TF-IDF)

The phrase "Tf-Idf" stands for "frequency-inverse document frequency," therefore Tf-Idf weight may also be a frequently employed weight for text mining and information extraction. This weight may be a numerical value used in a celebration or corpus to rate the significance of a word in a written piece.

The frequency with which a word appears in the document affects its importance, but the volume inside the body of the document is taken into consideration. Tf-Varying whether weighting is still the main method used by search engines to mark and rank documents in response to user queries.

The inclusion of the Tf-If for any query phrase is one of the sole ranking functions. The most complicated ranking functions are those that derive from this fundamental paradigm. Tf-If is frequently used to filter stop words in a variety of topic areas, including text and categorization.

A mathematical formula might be TF-IDF [12], which, among other things, corresponds to a term's value on a variety of records. The Reciprocal Text Frequency and Term Frequency combine two scales. The initial estimate suggests the following place by taking into account the occasions when the term is reported in the actual document. However, this could result in frustrating paradoxes, such the incredibly common "and." As a result, by weighing the duration of their case within an archival corps, the IDF equation offers greater motivation for people who occasionally resurface.

The tf-idf Vectorizer module from the sklearn library was used in this work to vectorize data [13].Additionally, the Vectorizer is in charge of assembling meaningless words, also known as NLP stop words, such as 'a,' 'a,.' 'in,' 'you,' 'them,' 'have,' or 'been,' with an unknown sense.

On the other hand, TF-IDF refers to a method of visiting a word in a massive archive. It attempts to attribute some of the desire to talk to the word itself. This is frequently employed in the mining

of raw materials. This weight can be a valid indicator of how significant a word is to a report in a corpus or collection.

Another technique for determining the subject of a piece of literature by counting the words it includes is the TF-IDF. As a result, the text appears less frequently. Terms were weighed using the TFIDF check value rather than the check frequency. In other words, the data set replaces word counts with TF-IDF ratings.

The first calculation made by TF-IDF is the "term frequency"—the number of times a term appears in a particular text. However, because words like "and" or "the" always act as complete documents, they must be consistently rejected. This is a part of the inverted document frequency. The importance of a term to differentiate one document from another decreases when it appears in more texts. Only the common and recognizable terms should be used as identifiers when traveling overseas. Every word's tf-idf significance might also take the shape of a structured arrangement that totals at least one.

In mathematically we can define the formula as:

$$W_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

tf I, j = Number of Occurrences of I in j
df i = Number of documents containing i

N = Total Number of documents

**Here we can define term frequency and inverse document frequency in a more specific way as follow:-**

**TF: Term Frequency,** What calculates the word count in a text? The term might be used significantly more frequently in longer texts than in shorter ones because each one differs in length. In terms of how to standardize by document length, the word frequency is also divided (for example, the entire number of words within the document)

# TF (t)

$$= \frac{Number\ of\ Times\ t\ appares\ in\ a\ document}{Total\ number\ of\ term\ in\ the\ document}$$

**IDF: Inverse Document Frequency,** Word value tests are what matter. When TF is being calculated, all terms are considered to be equally important. It is understood that words like

"is" or "is" may appear multiple times without denoting anything. Therefore, while measuring the proportion of rare expressions, we prefer to specify the frequent expressions.

# IDF (t)

$$= \log_e \frac{Total \ number \ of \ document}{Total \ number \ of \ document \ with \ term \ t \ in \ it}$$

# TF-IDF Vectors can be generated in light of any kind of input tokens such as words, characters, ngrams. The TF-IDF has three major levels as follow. Words NGram Character

Short Description of those Level in below.

- Words Level TF-IDF: Terms TF-IDF amount, each term in a matrix format is displayed as a TF-IDF value.
- N-gram Level TF-IDF: N termism mixture spoken at N-gram level. This matrix displays N-gram scores for TF-IDF
- Character Level TF-IDF: TF-IDF values of the n-grams character level in corpus represented in a matrix.

## 3.4 Model Generation with Algorithm

For this suggested application, eight different machine learning methods are used, and our programming language is Python 3.8.6. The Decision Tree (DT), Support Vector Machine (SVM), Random Forest (RF), K-Neighbors (KNN), Nave Bayes (NB), Passive Aggressive Classifier (PAC), Logistic Regression (LR), and Extreme Gradient Boosting (XGB) models and classification reports we made for the aforementioned data set have also been applied to assess how well our data matches the model. Using the practice data set, we updated our model for this machine learning technique, and after determining the output, we used the test data set to forecast our model. This action can apply to either algorithm or to both. Multiple classes can benefit from these algorithms, and distinct databases keep track of their characteristics and usage. The most popular ordering formulas, as we previously stated for Naive Bayes, are Logistic Regression

<center>

**CHAPTER 4**

# ALGORITHM

</center>

This section refers to the algorithms that were used during the classification process

## 4.1 Support Vector Machine

Another machine learning approach that may be seen and categorized into one of two classes in the results is a vector support machine, or SVM.

It is also possible that the Vector Machine Supports are a monitored and linear machine learning approach, which is frequently applied to solve classification problems. Additionally, a subset of SVRs known as Vector Regression Support use the same techniques to break down regression-related issues.

The Vector Machine Supports may also be a monitored and linear machine learning technique, which is widely used to address classification issues. The same methods are also used to deconstruct regression-related difficulties by a subset of SVRs called Vector Regression Support. Depending on which side of the hyperplane selection boundary the knowledge point lies, the appropriate class is given the knowledge point. The size of the features increases as the hyperplane becomes more complex. The direction and orientation of data points close to the hyperplane are controlled by what are known as vectors.

This was the project with the SVC module of the Scikit-reading library [12] that introduced the SVM classifier algorithm

**Algorithms for the Support Vector Machine work directly with the expression given below:-**

$$f(\text{w, b}) = \left[\frac{1}{n}\sum_{i=1}^{n}\max\left(0, 1 - y_i(w^T x_i - b)\right)\right]$$

## 4.2 Decision Tree (DT)

Decision trees (DT) are frequently employed in decision analysis and machine learning. A decision-making approach is a tree-like layout of options that includes accident outcomes, resource costs, and utility [35]. The internal nodes of the Tree Option represent the attribute state. Every internal node splits into branches that desire the result until it reaches the maximum point where it ceases to split and winds up to the blade nodes that reflect the range of the designated category. The DecisionTreeClassifier module of the sklearn library [13] utilized the DT algorithm.

The formula Among the supervised algorithm family is the decision tree. Unlike other supervised learning algorithms, the tree selection approach is frequently employed to address regression and classification problems. By learning fundamental judgments from prior data to forecast the category or value of the target variable, the Tree Option is used to create a model training that can be applied (training data). In decision trees, we begin by imagining a tree that can forecast a record's category mark. We contrast the values of the record attribute with those of the basic attribute. By following the branch in this manner, we go to the next level

.

Decision trees are often of two types:-
• Categorical Variable Decision Tree: Decision Tree which includes a categorical target of variable then it called a Categorical DT.
 • Continuous Variable Decision Tree: the choice Tree incorporates a continuous target of variable then it's called Continuous DT.


**For solving the choice tree problem for researchers mush using some criteria as follows**:-
 • Entropy
• Information gain
• Gini index
• Gain Ratio
 • Reduction in Variance
 • Chi-Square


**In Mathematically Entropy for "One" attribute is represented as:**

$$E(\mathrm{S}) = \sum_{i=1}^{c} p_i log_2 p_i$$

### 4.3 K-Nearest Neighbors

A classification and regression approach utilized is K-Nearest Neighbors (KNN). Equivalent data points in this algorithm's training cycle are regarded as being close to another
1. This proximity is determined using a variety of distance units, including Euclidean, Murkowski, Manhattan, etc. When a date cannot be categorised, the classifier conducts a voting task by majority, taking into account the proximity of the k most pertinent data points, and chooses the category mark that will be given at that moment.

The metrics described above are also a result of the significance of the k value when determining its value as necessary to increase its correctness [3]. To reach a conclusion, numerous tests are required for the methods used to find these values, which are commonly referred to as parameter tuning. The Scikit-learn library's NearestNeighbors module was used in this research to develop an algorithm from K-Nearest Neighbors [13].

Another straightforward technique is K nearest neighbors, which classifies newly discovered cases using a similarity metric and saves all existing examples (e.g., distance functions). Early in the 1970s, KNN was already being utilized as a non-parametric technique for statistical estimates and pattern identification. A majority of the case's neighbors vote to decide the case, and the case is then assigned to the most popular category among its closest K neighbors as defined by the distance function. If K is equal to 1, the case is just put in the group of its closest neighbor.

There are many distance functions already available to calculate the Neighbors algorithm for distance.

**Among all of them i've got used the Euclidean Distance formula as follows:-**

$$\text{Euclidean} = \sqrt{\sum_{i=1}^{k}(x_i - y_i)^2}$$

Also to be highlighted is the fact that these distance measurements are only applicable to the continuous variables we employed in our experiment. The Hamming distance must be employed when categorical variables are involved.

<div align="center">

**Hamming distance**

</div>

$$D_{\text{H}} = \sum_{i=1}^{k}|x_i - y_i|$$

## 4.4 Logistic Regression

When a variable is binary, the appropriate statistical method is logistic regression (binary). Like other regression analyses, logistic regression is a statistical analysis. The link between a single binary dependent variable and one or more individual nominal, ordinal, interval, and ratio-level variables can be explained and clarified using logistic regression. As a statistical regression, the statistical technique is often utilized for classification tasks (LR).The fundamental goal is to assign fresh observations a specific set of complexity marks backed by the concept of probability [12]. This is also accomplished by utilizing the prediction function, also referred to as the Logistic

function, which is called sigmoid, to transfer the output values of LR to the likelihoods. The decision to give a likelihood a special class mark will be determined by a parameter called a threshold that serves as a selection boundary.The predicted sting value is 0.5, however it can vary according to a variety of factors, including the metrics used for preparation, the necessity for tweaking, and the expected probability or class distribution. This project uses the LogisticRegression Sklearn package and a linear model to build the LR [13].

**The logistic regression estimates in mathematics that a multifunctional regression function is defined as**:

$$\textbf{\textit{Function (f)}} = \textbf{\textit{Log}}\left(\frac{P\,(V=1)}{1-(P=1)}\right) = \beta_0 + \beta_1 - x_{i_1} + \beta_2 - x_{i_2} + \underline{\quad\quad} + \beta_p - x_{i_m} \quad \text{........ [36]}$$

## 4.5 Random Forest

A community-based classification system called Random Forest (RF) has several decision-making bodies. The foundation of RF is the idea that any decision tree may predict the label of the output given the input, and can then attribute a wise crowd inspiration to the label that is most predicted. The results of the ensemble classification should be better than a random guess, but predictions from different decision-making bodies shouldn't be connected. Because trees actively try to correct for the many mistakes of each other, uncorrelated tree models are more accurate when integrated immediately as an ensemble model than any single model.The RF technique has been implemented using the RandomForestClassifier module in the Sklearn Ensemble Library [37]. Comparable to statistical regression employing a representational equation is logistic regression. In order to approximate the output value using weights or coefficient values (also known as the Athens letter beta), input values (x) are employed linearly (y). The primary distinction between a linear return and a numerical return is that the output value of a linear return is also binary (0 or 1) rather than numerical.

This implies that the likelihood still varies from 0 to 1. The probability of defaulting payment and not defaulting payment shall amount to 1 in the case of binary classification.

**Note: For binary classification and multi-class classifications, logistic regression may be used. An example of logistic regression is below:-**

$$y = e^{(b_0+b_{1*}x)}(1+e^{(b_0+b_{1*}x)}) \quad \dots\dots\dots \text{[38]}$$

## 4.6 Passive-Aggressive Classifier

A family of machine learning algorithms known as passive-aggressive algorithms should not be used.Beginner and even intermediate Machine Learning aficionados are familiar with it. Applications, however, would undoubtedly be extremely beneficial and efficient.

Note: This is also a high-level description of the algorithm that describes how it operates and when to use it. It's not going far into the mathematics of how it works.

In large-scale learning, passive-aggressive algorithms are frequently employed. It belongs to a small group of 'online learning algorithms. Because the computer file arrives sequentially in online machine learning methods, the machine learning model is gradually updated toward batch learning, where the complete training dataset is employed immediately. In situations where there are an unlimited quantity of data and it is computationally challenging to educate the complete data collection due to the sheer volume of data, this can also be very helpful.We'll only suppose that an algorithm for online learning will get a training example, update the classifier, and then discard the instance[39]. Finding fake news on a social media platform like Twitter, where fresh information is being posted every second, might be a very good example of this. The expertise required to continuously interpret Twitter data dynamically would be enormous, making an online learning algorithm the ideal tool. The fact that passive-aggressive algorithms don't require a learning rate makes them almost akin to Perceptron models. They do, however, contain a regularization parameter.

**Important parameters describe in below:**

C: This is the regularization parameter, and denotes the penalization the model will make on an incorrect prediction

**max_iter**: The maximum number of iterations the model makes over the training data.

**Tol:** The stopping criterion. If it is set to none, the model will stop when (loss > previous loss – to). By default, it is set to 1e-3.

- I have used the PAC algorithm to enforce the algorithm using the dataset to finding the best accuracy for the estimate solution for fake news detection. Mostly I have used it because:-
- It is fast! When I say it is fast I mean it takes ~1ms to train on 400 samples x 30 features each.  Accuracy varies. Greatly…..
- You all can achieve 100% accurate results on some datasets and something not.

My result of this algorithm given in result section

## 4.7 Ensemble Classifier

The custom classifiers were used in conjunction with the five custom classifiers listed above. The goal is to develop a voting classifier model that determines the weights to be applied to each classifier's prediction. Each training example is associated with the probability vector once the classifiers' chances are first kept for each training instance in the matrix. Then, a Logistic Regression model is given this vector matrix, which computes weights and generates a true (0) or false label result (1). Additionally, a voting classifier that uses a simple majority vote to select a winner among model predictions has been introduced. in contrast to the previously described ensemble model [10]

### 4.7.1 Extend Gradient Boosting Classifier

We will just train a model on our dataset using a daily machine learning model, a tree of choice, and then utilize it to make predictions. Even though we may have slightly changed the parameters or added more data, we are ultimately only utilizing one model. All models are trained and applied to our data independently while we build a series. On the other hand, boosting employs a more iterative methodology. Technically, it still uses an ensemble strategy, but it does it in a more intelligent way by combining multiple templates to execute the last template.

The GradientBoostingClassifier module from the Sklearn Ensemble Library has been used to create the XGB algorithm [13]. It was incorporated using this library. The benefit of this iterative method is that future model additions concentrate on fixing mistakes introduced by earlier models. Models are trained in isolation in the traditional ensemble process, which results in all models

making the same error [40]. In particular, the X-Gradient Boosting method trains new models to anticipate the leftovers of older models. Below is a diagram that I used to describe the method.



## 4.8 Multinomial Naive Bayes (MNB)

Due to the brevity and easy scaling of large scale tasks, Multinomial Naive Bayes (MNB) is also the classifier of choice for many text categorization problems. The output difference for contemporary discriminatory classifiers, however, is caused by robust data assumptions and does not affect classification accuracy. The MNB text categorization framework is examined in this research together with the most effective combination of common generation model adjustments. We examine direct search optimization utilizing random search techniques to improve the provided Meta parameter classifier.

In reality, many MNB implementations use modifications to disprove strong hypotheses, necessitating extra parameters known as meta parameters that are not part of the actual model. Class-conditional multinomial smoothing, a necessary modification of the MNB, may also serve as an illustration of this. The maximum likelihood estimates for words will result in zero estimates because to the information sparsity issue in very tongues, and a parameterized smoothing technique is required to successfully fix this. In conclusion, the MNB shape is utilized [41].

$$P\,(w,\,c) = P\,(w|c)\,p\,(c) \propto p\,(c) \prod_{n=1}^{N} p(n|c)^{W_n} \dots\dots\dots\dots \textbf{[41]}$$

# CHAPTER 5

# EVALUATION METRICS AND ANALYSIS

## 5.1 Performance parameters

A performance indicator might be a numerical declaration of both the representational labor and its results. Performance metrics are sponsored statistics that provide a clear picture of whether representation or action is succeeding in its aims and whether policy or organizational goals are being advanced.Precision, recall, F1-score, true negative rate, and false-positive rate accuracy were utilized as an evaluation matrix to gauge how well our suggested model performed.

### 5.1.1 Confusion Matrix

Confusion matrix is a table that, on occasion, does not accurately represent the results of a classification model (or "classifier") on a collection of test data that are believed to have accurate values.

Although the confusion matrix is often simple to comprehend, the related terms are constantly confusing. The information contained in the uncertainty matrix is also used to evaluate the classifier's performance.

There is a confusion matrix for the two-class problem in Table 1.

|  | Predicted Positive | Predicted Negative |
|---|---|---|
| Actual Positive | True Positive(TP) | False Negative(FN) |
| Actual Negative | False Positive(FP) | True Negative(TN) |

Table 1: Representation of the Confusion Matrix

### 5.1.2 Precision and Recall

Measurement of the capacity of the model to correctly classify the event of a positive class instance shall be calculated by recall. This is known as:-

**Recall** = **TP / (TP+FN)**

Beside, precision is:-

**Precision = TP / (TP+FP)**

### 5.1.3 F1-Score

F1 the score is the weighted average for Precision and Recall. Therefore, this score contains both false positive and false negatives in the estimate.

F1 = 2*(Recall * Precision) / (Recall + Precision)

### 5.1.4 True Negative Rate (TNR)

The true negative rate (TNR) is the proportion of samples tested negative using the test in question that are actually negative.

TNR =  TN / (FP+TN)

### 5.1.5 False Positive Rate (FPR)

**A False Positive Rate is an accuracy metric that can be calculated on a subclass of machine learning models.**

**FPR = FP / (FP+TN)**

**5.1.6 Accuracy**

Accuracy is a measure of the total of correctly defined samples taken from all samples. This is known as:-

Accuracy = ( TP+TN) / (TP+TN+FP+FN)

Now here can compare the classification matrix with our work as follow:-

• True Positive (TP): the amount of true positive examples is that the number of reports articles, correctly classified as fake;

• False Positive (FP): the amount of false-positive examples is that the number of reports articles incorrectly classified as fake;

• True Negative (TN): the quantity of true negative examples is that the number of stories articles, correctly classified as true;

• False Negative (FN): the amount of false-negative examples is that the number of reports articles incorrectly classified as true;

## 5.1.7 Result of a confusion matrix for selected Dataset

Using the feature extraction technique TF_IDF the Confusion Matrix for the selected Classification Model as follows:-

For CS-1 (Case Study)

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 4361 | 34 | 4332 | 63 | 4244 | 151 | 2824 | 1571 |
| 25 | 7304 | 28 | 7301 | 28 | 7301 | 226 | 7103 |

**Confusion Matrix for DT**　　**Confusion Matrix for SVC**　**Confusion Matrix for RF**　　　**Confusion Matrix for KNN**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 3122 | 1273 | 4339 | 56 | 4257 | 138 | 4374 | 21 |
| 23 | 7306 | 30 | 7299 | 53 | 7276 | 19 | 7310 |

**Confusion Matrix for MNB**　**Confusion Matrix for PAC**　　**Confusion Matrix for LR**　**Confusion Matrix for XGBC**

**Table 2: Representation of Confusion Matrix; DT: Decision Tree; SVM: Support Vector Machine; RF: Random Forest; KNN: k-Nearest Neighbors; MNB: Multinomial Naive Bayes; PAC: Passive Aggressive Classifier; LR: Logistic Regression; XGB: eXtreme Gradient Boosting**

**For CS-2 (Case Study)**

| 3320 | 23 | | 3282 | 61 | | 3204 | 139 | | 2087 | 1256 |
|------|------|---|------|------|---|------|------|---|------|------|
| 18 | 5432 | | 22 5 | 428 | | 17 | 5433 | | 158 | 5292 |

**Confusion Matrix for MNB**  **Confusion Matrix for PAC**  **Confusion Matrix for LR**  **Confusion Matrix for XGBC**

**Table 3: Representation of Confusion Matrix; DT: Decision Tree; SVM: Support Vector Machine; RF: Random Forest; KNN: k-Nearest Neighbors; MNB: Multinomial Naive Bayes; PAC: Passive Aggressive Classifier; LR: Logistic Regression; XGB: eXtreme Gradient Boosting**

**For CS-3 (Case Study)**

| 2195 | 11 | | 2153 | 53 | | 2135 | 71 | | 1336 | 870 |
|------|------|---|------|------|---|------|------|---|------|------|
| 10 | 3647 | | 22 | 3635 | | 13 | 3644 | | 94 | 3563 |

**Confusion Matrix for DT**  **Confusion Matrix for SVC**  **Confusion Matrix for RF**  **Confusion Matrix for KNN**

| 1449 | 757 | | 2165 | 41 | | 2108 | 98 | | 2199 | 7 |
|------|------|---|------|------|---|------|------|---|------|------|
| 7 | 3650 | | 26 | 3631 | | 28 | 3629 | | 7 | 3650 |

**Confusion Matrix for MNB**  **Confusion Matrix for PAC**  **Confusion Matrix for LR**  **Confusion Matrix for XGBC**

**Table 4: Representation of Confusion Matrix; DT: Decision Tree; SVM: Support Vector Machine; RF: Random Forest; KNN: k-Nearest Neighbors; MNB: Multinomial Naive Bayes; PAC: Passive Aggressive Classifier; LR: Logistic Regression; XGB: eXtreme Gradient Boosting**

# CHAPTER 6


# RESULTS AND DISCUSSION

## 6.1 Environments and Tools

In this section, the details of the environment and toolkit that were used for the implementation are briefly described below-:

## Software and Hardware configuration:

The implementation of proposed Framework was performed on Processor: AMD Ryzen 5 5500 3.6GHz-4.2GHz 6 Core 19MB Cache AM4 Socket Processor Installed RAM 16.00 (15.67 GB usable), system type 64-bit operating system, X, running under Windows 10 Pro operating system. The algorithm was in-house developed using Python-based Jupyter Notebook 2020 software. Apart from having used Matplotlib tools, Pandas tools, Sci-kit learn tools, and so on.

| ITEM | DETAILS |
|---|---|
| System Mode | Dextop |
| OS | Windows 10 Professional |
| Processor | AMD Ryzen 5 5500 3.6GHz-4.2GHz 6 Core 19MB Cache AM4 Socket Processor |
| CPU | |
| RAM | 16 GB |
| SSD | 1 TB |
| System Type | 64 |
| Tools | Jupyter Notebook, Pandas for data analysis and processing, Matplotlib for visualization, Scikit-learn, and Seaborn for advanced visualization. |

Table 5: Implementation details

## 6.2 Performance Comparison

The fake news dataset section 3.1 was initially used in a number of tests to test the output matrix of several machine learning algorithms. Matrixes of confusion Tables 2, 3, and 4 are additionally utilized for the particular machine learning classifier's performance assessment. parameters for performance evaluation 5.1.1 to 5.1.6. In this instance, our models have been successfully tested using a variety of output parameters. According to the arithmetic, calculations, analysis, and Page | 30 Daffodil International University inquiry, it was discovered that the classifier eXtreme Gradient Boosting had the maximum accuracy of 99.66 on Case Study-1 and also performed the best on Case Study -2 and Case Study -3. For our chosen dataset, machine learning-based classification outcomes are displayed as (Fig-11). In order to validate our classification results, we have also included parameters from these studies for measuring the various performance of any classifier, such as Precision, Recall, F1-Score, True Negative Rate, and False Negative Rate. Extreme Gradient Boosting (XGB) has been found to have a greater true negative rate and a lower false positive rate than all of our machine learning-based models. On the other hand, KNN has the largest false-positive rate whereas LR has a lower genuine negative rate (Table 2). With XGB as a classifier, we have been confirming our findings using other performance metrics like precision, recall, and F1-score. According to performed our investigation with machine learning-based models has found that performance decrease as the scale of data increases and after a certain period of data level performance increase gradually. As already said before in the abstract section where mentions that the dataset divided into three user cases of data.

**The Classification results for all machine learning models have been selected for CS-1, CS-2, and CS-3 (Section 3.1) as Illustrates in Fig-11.**



**SCORES AMONG ALL CLASSIFIERS**

| | DT | SVM | RF | KNN | MNB | PAC | LR | XGB |
|---|---|---|---|---|---|---|---|---|
| CS-1 | 99.49 | 99.22 | 98.47 | 84.67 | 88.95 | 99.27 | 98.37 | 99.66 |
| CS-2 | 99.53 | 99.06 | 98.23 | 83.92 | 87.49 | 99.09 | 98.12 | 99.82 |
| CS-3 | 99.64 | 98.72 | 98.57 | 83.56 | 86.97 | 98.86 | 97.85 | 99.76 |

CLASSIFIERS

©Daffodil International University

Figure 11: Different Machine Learning Accuracy Results; CS: Case Study

As has already been demonstrated, for the provided three-part dataset, the classification models employed for the inquiry in this publication perform best for XGB models and worse for KNN models. (See Section 3.1.) The important distinction is that the suggested text gradually gets shorter while the classification of DT gets bigger. On the other hand, when the proposed document steadily shrank, the performance of the SVM, KNN, MNB, PAC, and LR classifiers dropped.

| Classification Model | CS-1 | | | CS-2 | | | CS-3 | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1-Score | Precision | Recall | F1-Score | Precision | Recall | F1-Score |
| DT | 99.43 | 99.22 | 99.53 | 99.46 | 99.31 | 99.38 | 99.54 | 99.50 | 99.51 |
| SVM | 99.35 | 98.56 | 98.95 | 99.33 | 98.17 | 98.74 | 98.98 | 97.59 | 98.28 |
| RF | 99.34 | 96.56 | 97.93 | 99.47 | 95.84 | 97.62 | 99.39 | 96.78 | 98.06 |
| KNN | 92.59 | 64.25 | 75.85 | 92.96 | 62.42 | 74.68 | 93.42 | 60.56 | 73.48 |
| MNB | 99.26 | 71.03 | 82.80 | 99.42 | 67.48 | 80.19 | 99.51 | 65.68 | 79.13 |
| PAC | 99.31 | 98.72 | 99.01 | 99.51 | 98.02 | 98.75 | 98.81 | 98.14 | 98.23 |
| LR | 98.77 | 96.86 | 97.90 | 99.04 | 95.99 | 97.49 | 98.68 | 95.55 | 97.08 |
| XGB | 98.63 | 98.49 | 98.55 | 99.79 | 99.73 | 99.75 | 99.68 | 99.68 | 99.68 |

Table 6: Recall, Precision and F1-Score for Machine Learning algorithm; CS: Case Study; DT: Decision Tree; SVM: Support Vector Machine; RF: Random Forest; KNN: k-Nearest Neighbors; MNB: Multinomial Naive Bayes; PAC: Passive Aggressive Classifier; LR: Logistic Regression; XGB: eXtreme Gradient Boosting.

**In the meantime have seen before, using the formula from the classification matrix have found the highest precision, F1-score, and recall for the Decision Tree classification model of CS-1. For CS2 have been found highest recall, precision, and f1-score only for XGB classifier. For CS-3 given the same result 99.98 as the highest result for the precision, recall, and f1-score (Table 6)**

| Word Embedding Model | Classification Model | CS-1 Accuracy | CS-2 Accuracy | CS-3 Accuracy |
|---|---|---|---|---|
| Term Frequency–Inverse Document Frequency (TF-IDF) | DT | 99.49 | 99.53 | 99.64 |
| | SVM | 99.22 | 99.09 | 98.72 |
| | RF | 98.47 | 98.23 | 98.57 |
| | KNN | 84.67 | 83.92 | 83.56 |
| | MNB | 88.95 | 87.49 | 86.97 |
| | PAC | 99.27 | 99.09 | 98.86 |
| | LR | 98.37 | 98.12 | 97.85 |
| | XGB | 99.66 | 99.82 | 99.76 |

**Table 7: Accuracy of Word Embedding Model (TF-IDF) for Machine Learning algorithm.**

With three datasets (Section 3.1), the word embedding model TF-IDF in this experiment produced the greatest accuracy of 99.66, 99.82, and 99.76 (Table 7) for the XGB classifier among all classification models. In addition, KNN model accuracy has been determined to be the lowest. It is stated plainly that Boosting Technique will assist in determining the best estimate output solely of accurately detecting bogus news if our model is bound applying on the framework to do so. In addition, the performance of the DT classifier will be at its best as the dataset gets smaller. Several documents are dependent on the projected outcome for the DT classifier. As the dataset for the classifier, SVM, RF, and others, steadily shrinks, model performance for detecting false news will also decline.
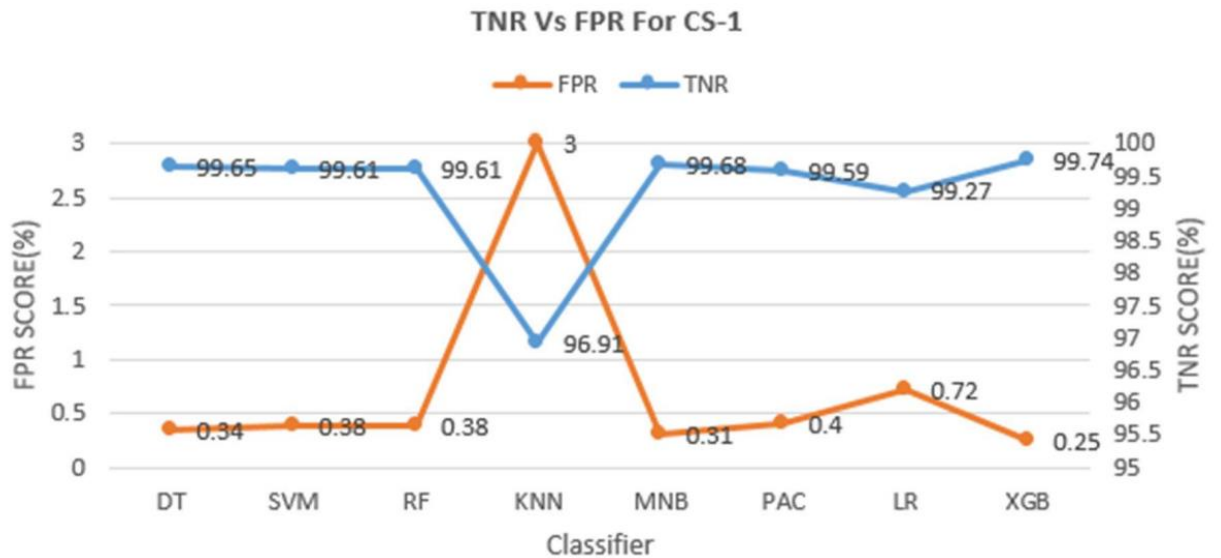


Figure 12: TNR vs FPR Score among all classifiers for CS-1; FPR: False Positive Rate; TNR: True Negative Rate
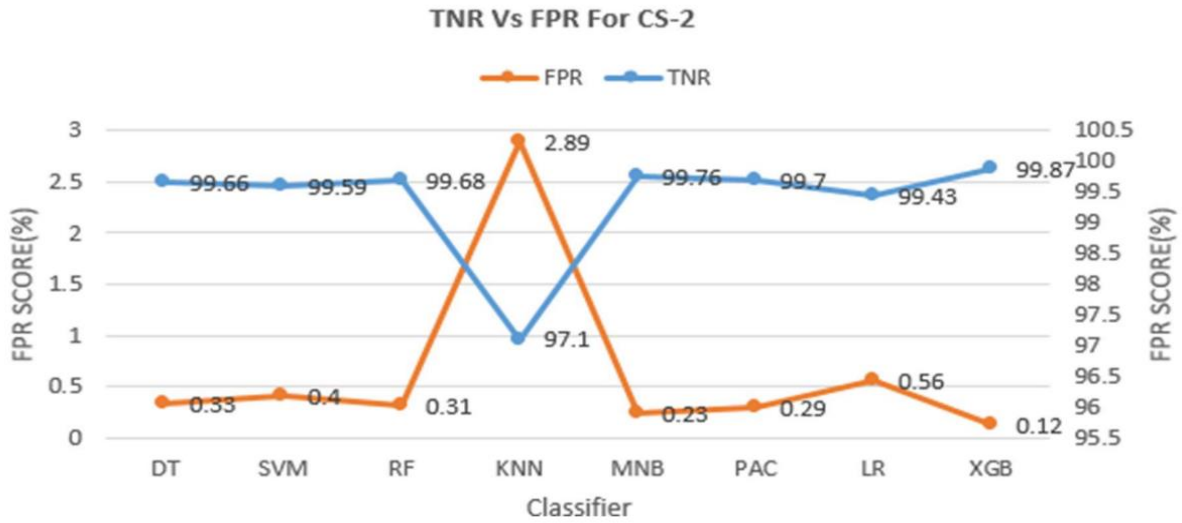
**TNR Vs FPR For CS-2**

Figure 13: TNR vs FPR Score among all classifiers for CS-2; FPR: False Positive Rate; TNR: True Negative Rate
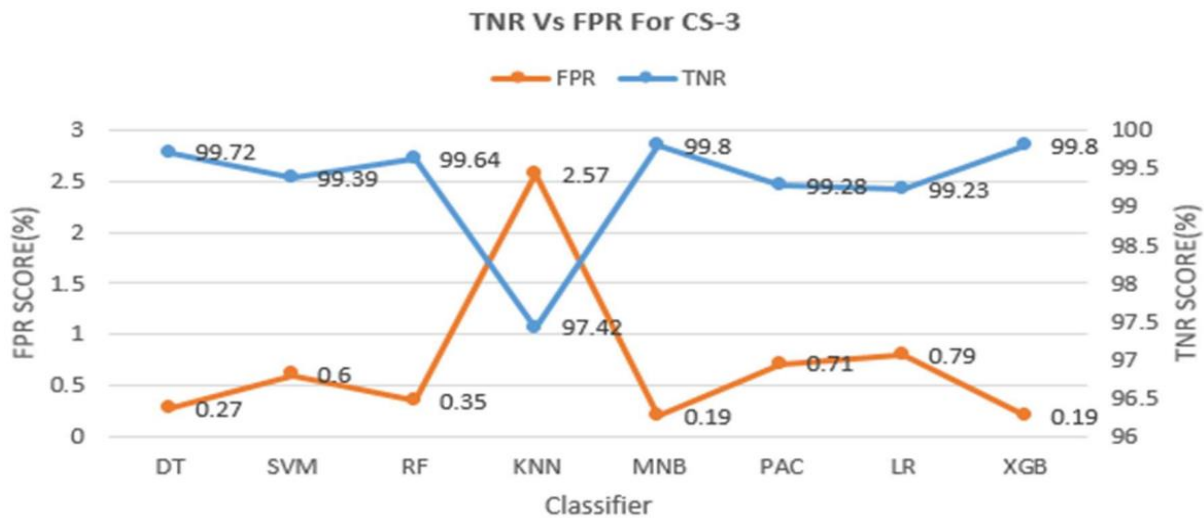


**TNR Vs FPR For CS-3**

Figure 14: TNR vs FPR Score among all classifiers for CS-3; FPR: False Positive Rate; TNR: True Negative Rate Since the total number of accurate negative predictions is divided by the total number of negatives, the specificity, or true negative rate (TNR), is calculated. The true negative rate is another name for it (TNR). The least specificity is 0.0, while the most specificity is 1.0. However, the rate of false-positive results (FPR) is calculated by dividing the total number of negative predictions by the number of inaccurate positive predictions. Simple false positive rates range from 0.0 to 1.0, with 1.0 being the worst. It can also be computed as 1 - TNR. As we can see from the aforementioned figures 4, 5, and 6, the XGB classifier has a maximum true negative rate and a minimum true negative rate. As a result, we have discovered the ideal model to identify false information for the XGB.

| Authors | Proposed Model | Environment | Testing Accuracy (%) |
|---|---|---|---|
| | | I | |
| Arvinder et al. (1) | XGB | Intel processor, core i7, DDR4 8GB Core 19MB Cache AM4 Socket Processor | **89%** |
| Dimitrios et al. (6) | CNN | MSI GeForce GTX 1630 VENTUS XS 4G OC 4GB GDDR6 Graphics Card | **89%** |
| Shaban et al. (7) | Hybrid Machine Crowd Approach | NA | **84%** |
| Rohit et al | . (8) FNDNet | NA | **98.36%** |
| OWN | FakeSpy | AMD Ryzen 5 5500 3.6GHz-4.2GHz 6 Core 19MB Cache AM4 Socket Processor | **98.99%** |

Table 8: Comparison-based Classification result using Kaggle Fake News Dataset

 Table 8 provides examples of several authors and their proposed models that have been used. according to their model. The setup of the computing environment for experimentation is provided below. The proposed Framework displays the most effective and comparative results (training). accuracy, testing accuracy, and a portable training model). This Framework is crucial for identifying bogus news. With the proposed Framework, we were able to obtain a 99.66% accuracy rate. Comparing this methodology to previous efforts, it produced better results with the real-world text-based fake news dataset. The above chart made it very evident that our suggested methodology is the most important framework to  identify false news, and we strongly encourage all studies to use it.

# CHAPTER7

# CONCLUSION AND FUTURE SCOPE

## 7.1 Conclusion

The deployment of eight TF-IDF machine learning algorithms to identify bogus news was demonstrated in this paper. We examined a computerized model that offers broad replies to data gathering and interpretative data for identifying fake news to verify the confirmation of information separated from the data collection.Based on the results of the aforementioned experiment and the coding of all machine learning classifier models, it was discovered that KNN had the lowest performance measurement and was the best classifier for identifying fake news.

Ensemble strategies will work best to address the major problems that have a fallout around the globe for detected fake news. Additionally, we may confirm our results using the confusion matrices and F1-score counts that we have constructed for each classifier model.The paper explains in detail how the XGB classifier can accurately identify fake news with 99.66% accuracy and 98.55% f1-score for the CS-1, 99.82% accuracy and 99.75% f1-score for the CS-2, and 99.76% accuracy and 99.76% f1-score for the CS-3 using a hot and novel set of features extracted from the heading and the text. We will suggest utilizing a decision tree classifier to identify bogus news for the smaller set of data.On the other hand, several ML classifiers have performed exceptionally well compared to the XGB classifier. Because of this, even though we used a large dataset, we chose the XGB models, saved this model, and used it for prediction. In order to identify bogus news, our suggested model successfully produces the predicted outcome. In the end, this significant issue might be successfully addressed using machine learning approaches. In this instance, the outcome strongly urges us to utilize our suggested methodology in the area of identifying fake news.

## 7.2 Future Scope

Other than deep learning and their model optimizers, various features may be applied to another feature extraction technique and a deep learning model with their optimizer and two or more features in the future for shaking. This may involve using a different two- or three-word embedding model to include a different model as well as additional linguistic features.

## REFERENCES

1. A. Balali, M. Asadpour and H. Faili, A Supervised Method to Predict the Popularity of News Articles, Computación y Sistemas 21 (2017), 703–716, ISSN 1405-5546.
2. M. Potthast, J. Kiesel, K. Reinartz, J. Bevendorff and B. Stein, A Stylometric Inquiry into Hyperpartisan and Fake News, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, 2018, pp. 231–240
3. N. Ruchansky, S. Seo and Y. Liu, Csi: A Hybrid Deep Model for Fake News Detection, in: Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, ACM, 2017, pp. 797–806
4. P. Krejzl, B. Hourová and J. Steinberger, Stance Detection in Online Discussions, arXiv preprint arXiv:1701.00504 (2017).
5. G. Sierra, Introducción a los corpus lingüísticos, Instituto de Ingeniería, UNAM: México, 2017, p. 210
6. W.Y. Wang, "Liar, Liar Pants on Fire": A New Benchmark Dataset for Fake News Detection, arXiv preprint arXiv:1705.00648 (2017).
7. T.H. Nazer, G. Xue, Y. Ji and H. Liu, Intelligent Disaster Response via Social Media Analysis A Survey, ACM SIGKDD Explorations Newsletter 19(1) (2017), 46–59
8. 3] J. Martínez, G. Bel Enguix and L. Torres Flores, Observations on Phonetic and Metrical Patterns in Spanish-language Proverbs, in: Proceedings of EUROPHRAS 2017, Tradulex, 2017, pp. 182–189
9. K. Shu, A. Sliva, S. Wang, J. Tang and H. Liu, Fake News Detection on Social Media: A Data Mining Perspective, ACM SIGKDD Explorations Newsletter 19(1) (2017), 22–36
10. B. Pang, L. Lee et al., Opinion mining and sentiment analysis, Foundations and Trends R in Information Retrieval 2(1–2) (2008), 1–135.
11. E. Stamatatos, A survey of modern authorship attribution methods, Journal of the American Society for information Science and Technology 60(3) (2009), 538–556.
12. Nielsen, R.K., et al., Navigating the 'infodemic': How people in six countries access and rate news and information about coronavirus. 2020: Reuters Institute
13. Reddy, H., et al., Text-mining-based fake news detection using ensemble methods. International Journal of Automation and Computing, 2020: p. 1-12
14. Star, T.D. A hazy picture appears. 2012.

15. Kalra, V. and R. Aggarwal. Importance of Text Data Preprocessing & Implementation in RapidMiner. in ICITKM. 2017.

16. Shah, F.P. and V. Patel. A review on feature selection and feature extraction for text classification. in 2016 international conference on wireless communications, signal processing and networking (WiSPNET). 2016. IEEE.

17.  . alokesh985, "Passive Aggressive Classifiers", geeksforgeeks, 2020, [Online]. Available: https://www.geeksforgeeks.org/passive-aggressive-classifiers/

# PLAGIARISM REPORT

## Turnitin Originality Report

Processed on: 15-Dec-2022 17:29 +06
ID: 1981930780
Word Count: 10057
Submitted: 1

**191-35-2685 By Hasibul Islam**

| Similarity Index | Similarity by Source | |
|---|---|---|
| **27%** | Internet Sources: | 24% |
| | Publications: | 21% |
| | Student Papers: | 5% |

---

8% match (Internet from 22-Apr-2022)
https://link.springer.com/chapter/10.1007/978-3-030-87954-9_15?code=89a1d995-3eb4-411b-a3b1-e48e3cb7fbde&error=cookies_not_supported

---

4% match ("Big Data Intelligence for Smart Applications", Springer Science and Business Media LLC, 2022)
"Big Data Intelligence for Smart Applications", Springer Science and Business Media LLC, 2022

---

2% match (Internet from 21-Nov-2022)
http://dspace.daffodilvarsity.edu.bd:8080/bitstream/handle/123456789/7544/173-35-2248%20%2823%25%29clearance.pdf?isAllowed=y&sequence=1

---

2% match (Internet from 04-Nov-2021)
https://coek.info/pdf-fndnet-a-deep-convolutional-neural-network-for-fake-news-detection-.html

---

1% match (Internet from 17-Oct-2022)
https://link.springer.com/chapter/10.1007/978-981-13-9942-8_40?code=d57e2768-0425-4a2c-9182-efa303929576&error=cookies_not_supported

---

1% match (Internet from 23-Jun-2022)
http://dspace.daffodilvarsity.edu.bd:8080/bitstream/handle/123456789/8238/181-35-2470%20%2819%25%29%20clearance.pdf?isAllowed=y&sequence=1

---

1% match (Internet from 26-Nov-2020)
https://www.geeksforgeeks.org/passive-aggressive-classifiers/

---

1% match (Shaban Shabani, Maria Sokhn. "Hybrid Machine-Crowd Approach for Fake News Detection", 2018 IEEE 4th International Conference on Collaboration and Internet Computing (CIC), 2018)
Shaban Shabani, Maria Sokhn. "Hybrid Machine-Crowd Approach for Fake News Detection", 2018 IEEE 4th International Conference on Collaboration and Internet Computing (CIC), 2018

---

1% match (Rohit Kumar Kaliyar, Anurag Goswami, Pratik Narang, Soumendu Sinha. "FNDNet – A deep convolutional neural network for fake news detection", Cognitive Systems Research, 2020)
Rohit Kumar Kaliyar, Anurag Goswami, Pratik Narang, Soumendu Sinha. "FNDNet – A deep convolutional neural network for fake news detection", Cognitive Systems Research, 2020

---

< 1% match (Internet from 26-Oct-2022)
http://dspace.daffodilvarsity.edu.bd:8080/bitstream/handle/123456789/8407/181-35-2376.pdf?isAllowed=y&sequence=1

---

< 1% match (Internet from 02-Nov-2022)
https://www.ryanscomputers.com/category/processor-amd

---

< 1% match (student papers from 29-Apr-2022)
Submitted to University of Hertfordshire on 2022-04-29

---

< 1% match (student papers from 23-Sep-2022)
Submitted to University of Hertfordshire on 2022-09-23

©Daffodil International University

# Account clearance