



GLOBAL JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY: C
SOFTWARE & DATA ENGINEERING
Volume 24 Issue 1 Version 1.0 Year 2024
Type: Double Blind Peer Reviewed International Research Journal
Publisher: Global Journals
Online ISSN: 0975-4172 & Print ISSN: 0975-4350

Innovative Approaches to Fake News Detection: A Data Mining Perspective

By Mustapha Ismail, Yayale Isihaka Muhammad
& Abdullahi Modibbo Abdullahi

Gombe State University

Abstract- Fake news becomes a major concern in the era of social media, as it can spread rapidly and has significant impacts on individuals and society. Society and individuals are negatively influenced both politically and socially by the widespread increase of fake news either generated by humans or machines. In the era of social networks such as Facebook, X (twitter) and WhatsApp, the quick rotation of fake news makes it challenging to evaluate its reliability promptly. Therefore, automated fake news detection tools have become a crucial requirement. To address the aforementioned issues, two data mining classification techniques were used as Extreme Gradient Boosting and Decision Tree with some python features. This study is designed to use Decision Tree and Extreme Gradient Boosting methods to develop an effective approach for detecting and classifying news as real or fake to obtain a reliable model performance. These models are trained on a labeled dataset consisting of both real and fake news.

Index Terms: *fake news, detection, data mining, social media, classification.*

GJCST-C Classification: *DDC Code: 006.312, 302.23, 384.3*



INNOVATIVE APPROACHES TO FAKE NEWS DETECTION DATA MINING PERSPECTIVE

Strictly as per the compliance and regulations of:



RESEARCH | DIVERSITY | ETHICS

© 2024. Mustapha Ismail, Yayale Isihaka Muhammad & Abdullahi Modibbo Abdullahi. This research/review article is distributed under the terms of the Attribution-NonCommercial-NoDerivatives 4.0 International (CC BYNCND 4.0). You must give appropriate credit to authors and reference this article if parts of the article are reproduced in any manner. Applicable licensing terms are at <https://creativecommons.org/licenses/by-nc-nd/4.0/>.

Innovative Approaches to Fake News Detection: A Data Mining Perspective

Mustapha Ismail ^α, Yayale Isihaka Muhammad ^σ & Abdullahi Modibbo Abdullahi ^ρ

Abstract- Fake news becomes a major concern in the era of social media, as it can spread rapidly and has significant impacts on individuals and society. Society and individuals are negatively influenced both politically and socially by the widespread increase of fake news either generated by humans or machines. In the era of social networks such as Facebook, X (twitter) and WhatsApp, the quick rotation of fake news makes it challenging to evaluate its reliability promptly. Therefore, automated fake news detection tools have become a crucial requirement. To address the aforementioned issues, two data mining classification techniques were used as Extreme Gradient Boosting and Decision Tree with some python features. This study is designed to use Decision Tree and Extreme Gradient Boosting methods to develop an effective approach for detecting and classifying news as real or fake to obtain a reliable model performance. These models are trained on a labeled dataset consisting of both real and fake news. The performance of the models was evaluated using standard evaluation metrics such as accuracy, precision, recall, and F1-score. The proposed approach achieved 100% accuracy in distinguishing between real and fake news. It revealed and highlighted the potential of utilizing data mining techniques to combat the spread of fake news and provide valuable insights for researchers and practitioners in the field of information confirmation/verification and media literacy. We hope to use a different dataset to test the proposed model.

Index Terms: fake news, detection, data mining, social media, classification.

I. INTRODUCTION

Have you heard or feel dilemma after confirming a particular news is false? False information is not new, however it has become a hot topic. Traditionally, we got our news from trusted sources, journalists and media outlets that are required to follow strict codes of practice. However, the internet has enabled a whole new way to publish, share and consume information and news with very little regulation or editorial standards. Many people now get news from social media sites and networks and often it can be difficult to tell whether stories are credible or not. Information overload and a general lack of understanding about how the internet works by people has also contributed to an increase in fake news or hoax stories (Yates, 2017). Also, (Tandoc, 2019) found that only

about 1% of examined Twitter account inspired 80% volume of sources of fake news.

Fake news refers to falsified news or propaganda disseminated through traditional media platforms such as print and television, as well as non-traditional media platforms such as social media [31]. The primary objective of disseminating such information is to deceive readers, harm a company's reputation, or profit from sensationalism (clickbait). It is widely regarded as one of the most serious threats to democracy, free speech, and social order [32]. Fake news is rapidly being disseminated through social media platforms such as twitter, Instagram, and Facebook, according to [33]. These platforms provide an avenue for the public to express their thoughts and opinions in an unfiltered and uncensored manner. Compared to conventional method views from media publishers' platforms, some news pieces hosted or shared on social media sites receive more views. According to researchers [32] who researched the speed with which fake news spreads on Twitter, tweets containing misleading information reach individuals six times faster than factual tweets. A few facts about fake news in the United States are as follows: 62% of Americans get their news from social media [34] making Fake news had a higher Facebook share than legitimate news [35].

Fake news detection is a subtask of text classification [38] and is often defined as the task of classifying news as real or fake. The term 'fake news' refers to the false or misleading information that appears as real news. It aims to deceive or mislead people. Fake news comes in many forms, such as clickbait (misleading headlines), disinformation (with malicious intention to mislead the public), misinformation (false information regardless of the motive behind), hoax, parody, satire, rumor, deceptive news, and other forms as discussed by [39].

Nowadays, information is easily accessible online, from articles by reliable news agencies to reports from independent reporters to extreme views published by unknown individuals. Moreover, social media platforms are becoming increasingly important in everyday life, where users can obtain the latest news and updates, share links to any information they want to spread, and post their own opinions. Such information may create difficulties for information consumers as they try to distinguish fake news from genuine news. The

Author α σ: Department of Computer Science, Gombe State University, Gombe, Nigeria.

Author ρ: Department of Computer and Information Science, Towson University, Towson Maryland, USA.

e-mail: aabdull2@students.towson.edu

wide spread of fake news on online social media has influenced public trust Knight Foundation, (2018), Naeem and Bhatti, (2020), etc. Under such severe circumstances, automatically detecting fake news has been an important countermeasure in practice. Based on a systematic review of recent literature published over the last five years, we synthesized different views dealing with fake news. We investigated machine learning (ML) applications to detect fake news, focusing on the characteristics of the different approaches and techniques, conceptual models for detecting fake news and the role of cognitive agents in this context as they have gained great popularity in the last few years. Data mining refers to extracting useful insights from large datasets, feature extraction is a technique that reduces raw data through extraction of most pertinent information [38], while ML algorithms are algorithms employed to learn from data and generalize to unobserved data [7].

This study was proposed to use a Decision Tree (DT) and Extreme Gradient Boost (XGBoost) machine learning algorithm to develop an accurate and reliable detection model that gives correct detection that is better than the one found in the published literature. Among other things the study seeks to address and show that a boosting algorithm outperforms other detection or prediction classifiers; design a framework for classification of news as fake or real using Boosting algorithms; describe what is fake news and how an individual can take preventive measures to avoid being contracted or being a victim; evaluate detection performances using evaluation metrics and compare with other results found in the published articles. The rest of the paper is organized as follows: Section II presents the existing relevant literature review; in section III the methodology employed was stated. Section IV provides the results obtained and discussion and finally, conclusion was drawn in section V.

II. LITERATURE REVIEW

The term fake news has recently become widespread. Even though there is no generally accepted definition of fake news, it continues to evolve day-by-day daily. Traditional fake news is generally defined as intentional behavior that harm a person or group, which could make it difficult for the victim to defend himself or herself. From the traditional definition of fake news, fake news could be explained as the use of information technology platforms especially social media to communicate wrong information about an individual or group, either intentional or otherwise. The use of news environment perception (NEP) to observe news environments for fake news detection on social media, designed popularity- and novelty-oriented perception modules to assist fake news detectors was proposed by (Sheng et al, 2022). Experiments on offline and online

data show the effectiveness of NEP in boosting the performance of existing models and drew insights on how NEP helps to interpret the contribution of macro and microenvironment in fake news detection [1]. T. SU (2022), proposed the use of a User Network Embedding Structure (UNES) model, which performs fake news classification on Twitter through the use of graph embedding to represent Twitter users' social network structure, Compared to the existing approach of using user networks with handcrafted features, UNES does not require any pre-annotated data (e.g., user type (individual users or publishers), users' stance, and if they have engaged with fake news before) and observed that using the user network embedding trained on a combined user network of two datasets is on par with or outperforms the user network embedding trained for the single experimental dataset on the MM-COVID and the SD datasets, respectively, which indicates the robustness of our proposed framework, FNDF. Thus, we showed that the three task models are all important components of our end-to-end fake news detection framework, and that the FNDF is robust when applied to news involving unseen users, if the user friendship network embedding is updated with the unseen users and their friends. Combining embedded entities with the language model results in as much as 177.6% increase in MAP on ranking check-worthy tweets, and a 92.9% increase in ranking check-worthy sentences [2]. In their study, Ali et al (2022), proposed and investigated several cutting-edge fake news detecting systems and associated problems. Methods for detecting and identifying false news, such as credibility-based, temporal-based, social context based, and content-based, were also thoroughly examined. Finally, the research investigates several datasets used to identify false news and proposed an algorithm [3]. Ahmad et al (2020), used ensemble techniques with various linguistic feature sets to classify news articles from multiple domains as true or fake. Ensemble techniques along with Linguistic Inquiry and Word Count (LIWC) feature set used in this research are the novelty of the proposed approach. There are numerous reputed websites that post legitimate news content, and a few other websites such as Politi-Fact and Snopes which are used for fact checking. In addition, there are open repositories which are maintained by researchers, accuracies of the techniques are: the accuracy achieved by each algorithm on the four considered datasets. It is evident that the maximum accuracy achieved on DS1 (ISOT Fake News Dataset) is 99%, achieved by random forest algorithm and Perez-LSVM. Linear SVM, multilayer perception, bagging classifiers, and boosting classifiers achieved an accuracy of 98%. The average accuracy attained by ensemble learners is 97.67% on DS1, whereas the corresponding average for individual learners is 95.25%. The absolute difference between

individual learners and ensemble learners is 2.42% which is not significant. Benchmark algorithms Wang-CNN and Wang-Bi-LSTM performed poorer than all other algorithms. On DS2, bagging classifier (decision trees) and XGBoost are the best performing algorithms, achieving an accuracy of 94%. Interestingly, linear SVM, random forest, and Perez-LSVM performed poorly on DS2. Individual learners reported an accuracy of 47.75%, whereas ensemble learners' accuracy is 81.5%. A similar trend is observed for DS3, where individual learners' accuracy is 80% whereas ensemble learners' accuracy is 93.5%. However, unlike DS2, the best performing algorithm on DS3 is Perez-LSVM which achieved an accuracy of 96%. On DS4 (DS1, DS2, and DS3 combined), the best performing algorithm is random forest (91% accuracy). On average, individual learners achieved an accuracy of 85%, whereas ensemble learners achieved an accuracy of 88.16 %.) Worst performing algorithm is Wang-Bi-LSTM which achieved an accuracy of 62% [4].

Althabiti et al (2022) examined an English dataset labelled as whether a particular article is 'true', 'false', 'partially false' and 'other', investigated four ML algorithms and pre-trained transformers to solve this multi-classification problem and attempted to use an external dataset from Kaggle to help improve the model. However, the additional dataset did not increase the performance, even though we used a different number of samples in each attempt. Finally, their findings from over 30 experiments show that the BERT model outperforms other models. The obtained testing results on the leader board indicate that we got an F1 of around 0.305, which slightly differs from the highest participant's score with only about 0.03. Future work recommended finding an additional dataset with a similar format may help improve the model. Also, using an ensemble method, which considers both rule-based and deep learning methods, could significantly enhance the proposed system [5]. The study by (Johnson1 et al, 2021), used random forest and decision tree algorithms on a dataset containing both fake and real news to do classification. The software used for the experiment was WEKA and the result generated showed that random forest correctly classified instance is 100% and incorrectly classified instance is 0% while the decision tree correctly classified instance is 93.6364% and incorrectly classified instance is 6.3636%. The results are a proof that random forest algorithm is a better classification tool as compared to decision tree. The results obtained show that Random Forest is a better classification tool with correctly classified instance of 100% and incorrectly classified instance of 0% as compared to the decision tree with correctly classified instance of 93.6364% and incorrectly classified instance of 6.3636%. It is recommended that future studies be carried out in the area of fake news prevention so that

fake news after being detected can be blocked from gaining access into the society. They used a classification report and confusion matrix to assess their model during the validation phase [6]. The work by (Sharma et al, 2020) aimed to perform binary classification of various news articles available online with the help of concepts pertaining to Artificial Intelligence, Natural Language Processing and ML. They also aimed to provide the user with the ability to classify the news as fake or real and also check the authenticity of the website publishing the news, various NLP and ML Techniques have to be used. The model is trained using an appropriate dataset and performance evaluation is also done using various performance measures. The best model, i.e. the model with highest accuracy is used to classify the news headlines or articles. As evident above for static search, our best model came out to be Logistic Regression with an accuracy of 65%. Hence they used grid search parameter optimization to increase the performance of logistic regression which then gave us the accuracy of 75%. As a result, they can say that if a user feed a particular news article or its headline in our model, there are 75% chances that it will be classified to its true nature. The user can check the news article or keywords online; he can also check the authenticity of the website. The accuracy for dynamic system is 93% and it increases with every iteration. We intent to build our own dataset which will be kept up to date according to the latest news [7]. E. K. Qalaja et al (2022), the authors employed supervised ML techniques on our newly developed dataset. Specifically, the proposed system categorizes fake news related to COVID-19 extracted from the Twitter platform using four ML-based models, including decision tree, Naïve Bayes (NB), artificial neural network (ANN), and k-nearest neighbors (KNN) classifiers), our experimental evaluation reported that DT based detection model had achieved the highest detection performance scoring 99.0%, 96.0%, 98.0%, and 90.0% in ACC, FSC, AUC, and MCC, respectively. The second set of experiments employs the small dataset (i.e., 700 tweets); their experimental evaluation reported that DT based detection model had achieved the highest detection performance scoring 89.5%, 89.5%, 93.0%, and 80.0% in accuracy, f1-score, area under the curve, and MCC, respectively. The results obtained for all experiments have been generated for the best-selected features [8].

Shu et al (2017), proposed the use of TriFN to detect fake news on social media where focused on using news contents and social contexts. For news content-based approaches, features are extracted as linguistic-based and visual-based. Linguistic-based features aim to capture specific writing styles and sensational headlines that commonly occur in fake news content Potthast et al. (2017), Afroz, Brennan, and Greenstadt (2018). For social context-based

approaches, the features include user-based, post-based and network-based. User-based features from user profiles to measure their characteristics and credibility (Castillo et al, 2011) and (Kwon et al, 2013). Finally came out with the “accuracy of 80%”. [9]. In their study (Unirio et al, 2019), applied neural network using WEIBO dataset to detect fake news on social media achieving the degree of accuracy seventy five percent [10]. Orellana et al (2018), the authors proposed the use of ML, text analytics and network models – to understand the factors underlying audience attention and news dissemination on social media (e.g., effects of popularity, type of day) and also provide new tools/guidelines for journalists to better disseminate their news via these social media [11]. According to (Bondielli et al, 2019), the use of tree like network using Breath First Search (BFS) strategy to analyze and summarize the approaches for source detection of rumor and misinformation in social network and provides an intense research contribution for further exploration of source detection of rumor in a social network [12].

According to (Shelke et al, 2019), the use of Data mining, ML, Classification application with automated fact-checking applications developed to tackle the need for automation and scalability and came out with the accuracy performance of classification models 88.2%[13]. The study published by (Nyow et al, 2019), proposed the use of Artificial Intelligence, Natural Language Processing and ML to provide the user with the ability to classify the news as fake or real and also check the authenticity of the website publishing the news with multiple models trained and also some pre-trained model extracted from Felipe Adachi. The accuracy of the model is around 95% for the entire self-made model and 97% for this pre-trained model [14]. Aphiwongsophon et al (2018), explored the used of ML techniques to detect fake news by using four popular methods in the experiments: - Naïve Bayes, neural network, SVM and the normalization method for cleaning data before using the ML method to classify data. The result shows that the Naïve Bayes used to detect fake news has accuracy. Two other more advanced methods which are neural network and SVM achieved the accuracy of 99.90% [15]. Rubin et al (2016), proposed the use of satire method, satire is a type of deception that purposely incorporates cues revealing its own deceptiveness, the deception detection was quite challenging. However, the method was able to integrate word level features using an established ML approach in text classification and SVM. The style-based deception detection method reaches relatively high accuracy rates of 90%, precision of 84% and recall of 87% [16]. Ray et al, (2017), considered the use of naïve Bayes classifier to detect fake news by Naive Bayes. This method has performed as a software framework and experimented it

with various records from the Face book, etc., resulting in an accuracy of 74% [17]. The paper neglected the punctuation errors, resulting in poor accuracy [18].

Gil, P (2019), the estimated various ML algorithms and made the researches on the percentage of the prediction. The accuracy of various predictive patterns included bounded decision trees, gradient enhancement, and SVM were assorted. The patterns are estimated based on an unreliable probability threshold with 85-91% accuracy [19]. Tandoc et al, (2017), utilized the Naive Bayes classifier, discussed how to implement fake news discovery to different social media sites. They used Face book, Twitter and other social media applications as a data source for news. Accuracy is very low because the information on this site is not 100% credible [20]. Sharma et al (2019), presented feedback-based approaches for fake news detection. In content-based approaches, the text of an article is regarded as the primary source of information. However, rich secondary information in the form of user responses and comments on articles and patterns of news propagation through social media can likely be more informative than article contents that are crafted to avoid detection. These secondary information sources form the basis of the works discussed [21]. Devi et al (2019), proposed the use of text processing and Naïve Bayes for training model and analyzed detection of fake news which is now prevalent in social media platforms and websites, used Therefore by using ML techniques and concluded that any news from large or small dataset can be classified as fake or not fake with previous data set values in less time which helped the user to believe in particular news that appears on social media or other sources [22].

Kesarwani et al (2020), proposed the use of a simple approach for detecting fake news on social media with the help of K-Nearest Neighbor classifier and achieved a classification accuracy of this model approximate 79% tested against Face book news posts dataset [24]. Khanam et al (2021), the authors proposed the use of six algorithms used for the detection are as: XGboost, Random Forests, Naive Bayes, K-Nearest Neighbors (KNN), Decision Tree, and SVM. The confusion matrix is automatically obtained by Python code using the cognitive learning library when running the algorithm code in Anaconda platform. Three common methods are utilized through their researches Naïve Bayes, Neural Network and SVM. Naïve Bayes has an accuracy of 96.08% for detecting fake messages. The neural network and the support vector machine (SVM) reached an accuracy of 99.90%. The scope of this paper is to cover the political news data, of a dataset known as Liar-dataset, it is a New Benchmark Dataset for Fake News Detection and labeled by fake or trust news. We have performed analysis on "Liar" dataset. The results of the analysis of the datasets using

the six algorithms have been depicted using the confusion matrix [25]. The current methods are reviewed for detection and classification of fake news using different supervised learning algorithms and a few unsupervised learning.

III. MATERIALS AND METHOD

Several classification algorithms can be used to classify whether given news is real or fake. But for this study, since we want to make a thorough detection, two different ML classification algorithms were chosen based on the published articles reviewed. These include eXtreme Gradient Boosting and Decision tree (DT). The eXtreme Gradient Boosting was employed as it has been optimized to increase GMB's speed and prediction performance; it is scalable and integrated into a different

platform. It is faster than other algorithms due to less resource usage. It has a new tree learning algorithm to handle fewer data while Decision Tree (DT) was used because it minimizes the chance of missing crucial information or taking the wrong steps. Which could lead to unnecessary escalation. Moreover, it equips frontline agents with the necessary knowledge to handle a range of inquiries confidently, mitigating the need to involve supervisors or specialized teams.

Model built was subjected on the training data to learn from it and evaluated on the testing data. The results obtained is evaluated on performance evaluation metrics for further determination and investigation of best performing model for fake and real detection. Figure 1 Demonstrates the Framework of the Study.

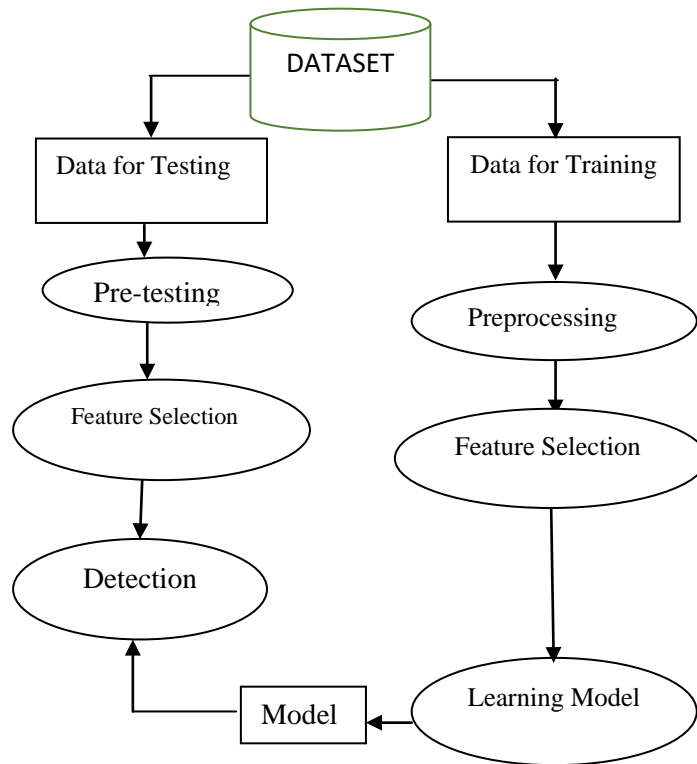


Figure 1: Framework of the Study

Experimental Setup

First of all, it was ensured that all necessary programs, tools, and techniques that would be needed to perform this experiment were downloaded and installed to obtain a good result. The python libraries used for various operations and functions are also installed. These include NumPy which is used for numerical python; pandas are used for data loading and analysis is also acquired and installed to set the environment ready. Jupyter notebook was used because it presents codes and data very well. Scikit-learn python machine learning library is used for designing, building; the system and software used for

this experiment is Windows 10, Python 3.7 and colab notebook were used to run the experiment.

Data Collection

Data collection is a critical step in the research process, as the quality and accuracy of the data collected impacts the validity and reliability of the findings. It is important to ensure that the data collection process is well-designed, carefully planned, and effectively executed to minimize errors and biases. The data is Downloaded in Comma Separated Value (Csv) Format.



Data Description

Data description is important because it helps to understand the dataset and its properties, which can guide the choice of appropriate techniques and methods for analyzing the data. It also helps to identify potential issues or problems in the data, such as outliers, missing data, or measurement errors.

Data Preparation

After importing the necessary supporting programs, files, and dataset for the research work. Data preparation is paramount which is used to set the data ready to go for a machine learning project. The data collected was loaded into the google colab notebook using a panda's command read data as follows.

```
#Import the data
from google.colab import drive
drive.mount('/mntDrive')
```

It is common knowledge that most data collected or downloaded must be prepared or preprocessed to make it fit the proposed model to obtain an accurate and even dependable result. Therefore, the data obtained has undergone data preprocessing, feature extraction and feature engineering as briefly explained.

The goal of data preprocessing is to improve the quality of the data, remove any inconsistencies or errors, and make it easier to work with. Data preprocessing is a crucial step in the data analysis process, as it has a significant impact on the accuracy and validity of the results. Properly preprocessed data helps to improve the performance of machine learning models and other analytical tools, leading to better insights and decisions preprocessing: In any Machine Learning process, data preprocessing is the step in which the data gets transformed, or encoded, to bring it to such a state that now the machine can easily parse it. In other words, the features of the data can now be easily interpreted by the algorithm. In this fake news detection, preprocessing is the major thing that should be done. Firstly, as the dataset is collected from kaggle.com. Therefore, unnecessary pieces of information were removed, converted to lower case, removed punctuation, symbols and stop words and so on.

a) Dataset and Data Preprocessing

Dataset was collected from a popular ML repository called kaggle with 44919 rows and 6 columns. It is the one of the largest community of data scientist in the world. Pre-processing refers to the transformations that were applied to data before feeding it to the ML algorithm. Data Preprocessing is a technique that is used to convert the raw data into a clean data set. Some of the important data preprocessing techniques used in the study was data

cleaning, dimensionality reduction, and feature engineering.

This is an essential phase that is used to enhance the quality of data to promote the extraction of meaningful insights. To also ensure that too much noise is minimized in the dataset to avoid over-fitting or under-fitting the proposed designed model. Improve the computational efficiency and accuracy of the model performance. To prepare the dataset for appropriate prediction, overfitting was avoided by training the model with sufficient number of rows and columns while under fitting was avoided by amping up model complexity, down regularization and data collection.

b) Feature Extraction

The feature extraction techniques was used in this research to reduce dataset over fitting the prediction model, improve prediction accuracy and reduce model training time. The dataset features that would be used to train the ML models have a great influence on the performance of the algorithm. Irrelevant, inappropriate or partially relevant features can undesirably influence model performance. Having unrelated features in the data can decrease the accuracy of the models, especially linear algorithms like linear and logistic regression. This feature extraction step is a process of dimensionality reduction process by which an initial set of data is reduced by identifying only the most relevant key features from the dataset that affects the detection machine learning model. In this study, data mining classification techniques were employed (DT and XBootst Algorithms) due to; DT is being faster than other algorithms due to less resource usage. It has a new tree learning algorithm to handle fewer data. Parallel and distributed computing accelerate learning, allowing for faster model discovery. Its prediction success is quite high while Gradient boosting decision trees are relatively easy to implement. Many include support for handling categorical features, don't require data preprocessing and streamline the process of handling missing data. Feature extraction was used to extract necessary relevant features for fake news detection and classification model.

Feature Selection

This feature is also used primarily to improve or enhance model detection performance accuracy. It is a technique of machine learning that leverages data to create new variables that are not in the training dataset. It generally produces new features for both types of machine learning projects, supervised and unsupervised learning. It is used for simplifying and increasing the speed of dataset transformation and manipulation aside from improving precision, recall, F1-score, and accuracy of model performance. Developing machine learning classifiers like extreme gradient boost and decision tree are also set in place.

Filter Method: Filter feature selection method and intrinsic techniques were employed to evaluate the relationship between each input variable and the target variable, and these scores are used as the basis to choose (filter) those input variables that were used in the model meanwhile intrinsic Algorithms that perform automatic feature selection during training (Decision Trees).

Model Application

There are a lot of already developed machine learning algorithms that has been rightly used and applied them directly for detection and classification purposes but for this study work, the extreme gradient boosting and decision tree algorithm were developed to fit the proposed model for this work. Hence, the modified extreme gradient boosting and decision tree algorithm models have been applied to the already preprocessed or prepared dataset for the detection and classification of news as fake or real.

Train-Test Split

Machine learning classification algorithms were used in training the model because the accuracy of the machine learning model mostly depends on the model training on the dataset. After the model has been developed and trained then testing is necessary to measure how accurate the detection performance of the model is.

At this stage, the dataset was divided into two sets, seventy-five and twenty-five percent; the first one for training and the former for testing, this was done to evaluate the model performance on the dataset that is not known to the model.

A Learning Model

It is a program that can find patterns or make decisions from a previously unseen dataset. For example, in natural language processing, machine learning models can parse and correctly recognize the intent behind previously unheard sentences or combinations of words.

Model

A model an informative representation of an object, example, pattern, exemplar, ideal mean someone or something set before one for guidance or imitation. Model applies to something taken or proposed as worthy of imitation.

It was observed that the actual parameters and their impact may vary depending on the specific implementation and version of the algorithms used. Additionally, the choice of Decision Trees and XGBoost in this study depends on the nature of the problem, the dataset size, and the desired balance between interpretability and predictive performance. Table 1 summarizes the comparisons of model parameters.

c) *Fine Tuning*

Fine-tuning typically refers to the process of taking a pre-trained ML model and training it further on a specific task or dataset to improve its performance. This is to take advantage of the knowledge and information the model has already learned from the large amount of data in the pre-training process. It allows you to transfer pre-existing knowledge to a new task where one can continue to train and adapt the model to improve its accuracy and efficiency.

The fine-tuning techniques consist of the following steps:

1. Loading the pre-trained model.
2. Freeze most if not all layers in the model to prevent them from further training.
3. Swap final layer or layers of the model with a new one that are specific to the task.
4. Train the model on dataset using a lower learning rate than in the pre-trained phase.
5. Evaluate performance of the fine-tuned model and adjust the hyper-parameters as necessary.



Table 1: Summary and Comparison of the Model Parameters

Parameter	Decision Trees	XGBoost
Learning Algorithm	Greedy recursive partitioning	Gradient boosting
Ensemble Method	Not an ensemble method	Boosting ensemble method
Regularization	Prone to overfitting	Includes regularization (L1, L2 penalties)
Handling Missing Values	Not naturally handled	Can handle missing values natively
Feature Importance	Provides feature importance	Provides feature importance
Parallel Processing	Generally, not parallelizable	Can be parallelized
Speed	Can be slower for large datasets	Faster due to parallelization and optimization
Robustness	Sensitive to noise and outliers	More robust due to ensemble and regularization
Hyperparameter Tuning	Fewer hyperparameters to tune	More hyperparameters to tune
Memory Usage	Lower memory usage	Higher memory usage
Suitable for Large Datasets	Limited scalability	Well-suited for large datasets

d) Prediction Tools

The python programming language would be throughout this study. Sci-kit learn libraries would be employed to help experiment. Python has a huge set of libraries and extensions, which are specifically designed for prediction models. Sci-kit learn is one of the best sources for ML algorithms <https://scikit-learn.org> where nearly all types of ML algorithms are readily available, easy, and quick evaluation of ML algorithms is possible. Numpy and Pandas will be used to deal with the data. For debugging and its ability to present code nicely Jupyter Notebooks was used.

e) Performance Evaluation Technique

The generality of the training datasets is the major goal of building a prediction model using ML techniques. ML models should be able to perform pretty well on real data. The dataset will be divided into two categories; training data and testing data. Training data will be used to train ML classifiers whereas testing data to test ML classifiers.

f) Evaluation Metrics

These are tools that are used to measure the effectiveness of the proposed model, to determine whether the built model can accurately make the required prediction. Many of these tools are in existence but the most commonly used for future predictions are; accuracy, precision, recall, and f1-score. They are calculated as follows:

g) Accuracy, Precision, Recall, And F-Score

These evaluation metrics were used to evaluate fake news detection models in (Poddar & D, 2019), (Ahmad et al., 2020), and (Ghafari et al., 2020). They are calculated as follows:

Accuracy is simply defined as the measure of the ratio of all testing samples which is classified as correct.

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn}$$

Precision means the ratio of relevant classified samples among the total retrieved samples.

$$Precision = \frac{tp}{tp + fp}$$

Recall is defined as the ratio of relevant classified samples among the total amount of relevant samples.

$$Recall = \frac{tp}{tp + fn}$$

F1-Score is the harmonic average of the precision and recall.

$$F1 - Score = 2 \left(\frac{precision * recall}{precision + recall} \right)$$

Where:

TP = True Positive, TN = True Negative, FP = False Positive and FN = False Negative values.



IV. RESULT AND DISCUSSION

Figure 1, presents the chart showing the number of real news (1) and fake news (0) available in the dataset. The fake news data has slightly outperformed the real news data in accuracy as shown in figure 2.

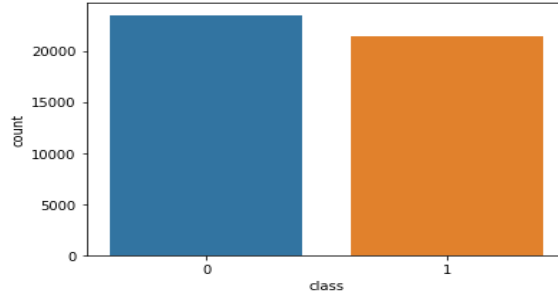


Figure 2: Count of Fake and Real News

A word cloud is a graphical representation of textual data that display the most frequently occurring words in a given text or dataset. The words are usually displayed in a randomized way as shown in figures 3 and 4, with their font size and color determined by their frequency of occurrence. Word clouds are often used to provide a quick visual summary of a large text or

dataset, highlighting the most important or relevant terms. Figure 3 and 4 displayed the word cloud of real news and word cloud of fake news respectively. The word with the highest frequency on figure 3 is “said” while on figure 4 is “trump”. The word “trump” appears frequently because the USA elections was trending then and made it appears much in the dataset used.

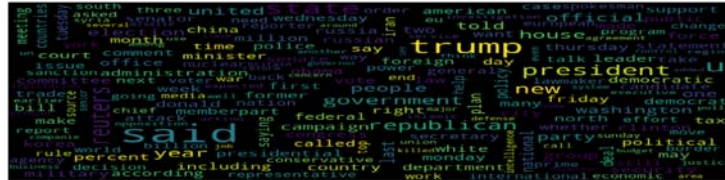


Figure 3: Word Cloud of Real News



Figure 4: Word Cloud of Fake News

Figure 5 displayed the bar chart of words frequency, it refers to the number of times a word appears in a given text or dataset. It is a measure of the importance or relevance of a word within the context of the text. The bar chart just as the word cloud indicated that the word said appeared the highest number of times in the data.

positive prediction and correctly made 5323 true negative prediction. It however, failed to make 14 false positive predictions and 32 false negative predictions. Hence, the decision tree algorithm performed very good prediction.

Confusion metrics are useful for evaluating the performance of classification algorithms, as they provide a more detailed understanding of the accuracy and errors of the model. From the confusion matrix, various performance metrics such as accuracy, precision, recall, and F1 score can be calculated. Figure 6 and figure 7 are confusion matrices of results obtained from decision tree classification and extreme gradient boosting classifier.

As presented in the confusion matrix figure 6, the decision tree classifier has correctly made 5861 true

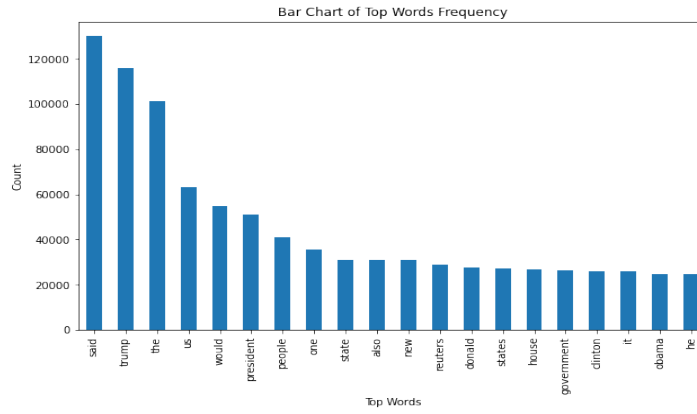


Figure 5: Bar Chart of Top Words Frequency

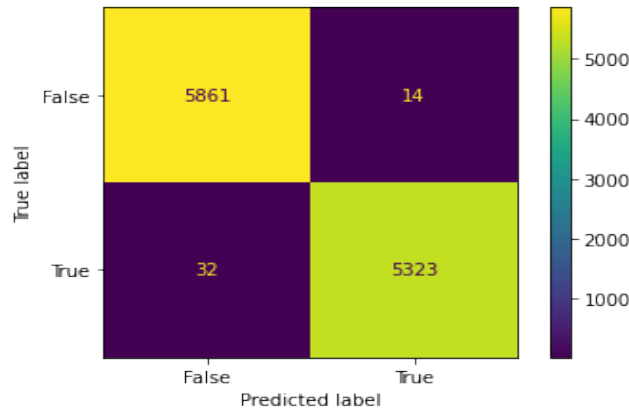


Figure 6: Confusion Matrix of Decision Tree Model

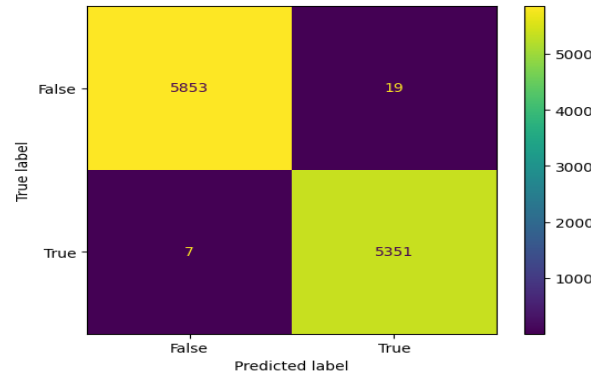


Figure 7: Confusion Matrix of Xgboost Model

As presented in the confusion matrix in figure 6 above, the extreme gradient boosting model has correctly made 5853 true positive prediction and also correctly 5351 true negative prediction. It only makes 19

false positive prediction and 7 false negative prediction. Therefore, this model has perform extremely well based on the result presented.

Table 2: Evaluation Results

	Model	Accuracy	Precision	Recall	F1-Score
Proposed	XGBoost	1.00	1.00	1.00	1.00
	DT	0.99	0.99	1.00	1.00
Existing	KNN	0.89	0.90	0.88	0.89
	DT	0.89	0.90	0.85	0.88

The simplest intuitive performance metric is accuracy, which is the ratio of properly predicted observations to all observations. Figure 7 showed the comparison of the entire various ML model used in the study as also shown in Table 2. As displayed in figure 7, accuracy, figure 8, precision, figure 9, recall and figure 10, f1-score. The extreme gradient boosting model have the highest prediction performance for all evaluation metrics which is also higher than the existing result

published by [43]. Therefore, for efficient prediction performance and decision-making, a precise forecast and classification of fake and real news is highly required. XGBoost is a powerful ML algorithm that had been employed in various applications and fake news detection with several contributing factors like gradient boosting, handling high dimensional data, regularization etc [29]. which possibly made it perform better that DT and KNN.

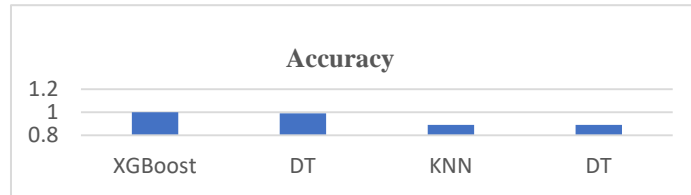


Figure 7: Accuracy of the Prediction Model

The ratio of accurately predicted positive observations to total expected positive observations is known as precision. As indicated in figure 8, the

XGBoost has the highest performance followed by decision tree.

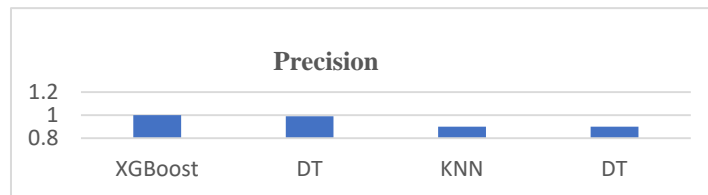


Figure 8: Precision of the Models

Figure 9 present the recall which is defined as the proportion of accurately predicted positive observations to all observations in the class. Again, the

proposed model that is extreme gradient boosting model has performed better than the existing model.

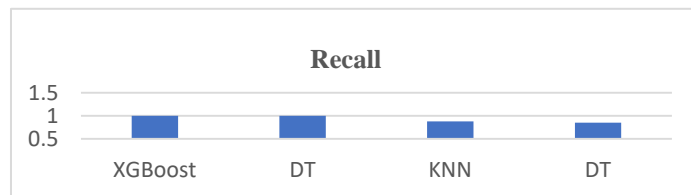


Figure 9: Recall Performance of the Models

Figure 10 shows the f1-score of the performance model. It considers both false positives and false negatives. Although it is not as intuitive as accuracy, F1-score is frequently useful than accuracy, especially if the class distribution is unequal. Our

proposed model also outperformed highly incredible compared to the existing model. From the results in table 2 F1 score and a recall of 1 respectively shows an excellent performance.

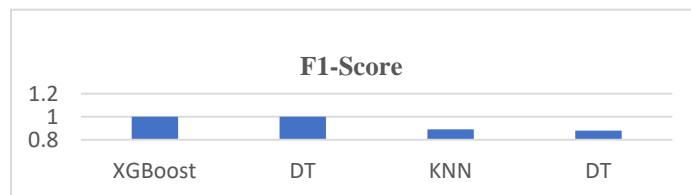


Figure 10: F1-Score Performance of the Models

The confusion matrix result presented in table 2 is a useful tool for understanding the performance of a classification model. It also allowed us to calculate several metrics that can be used to evaluate the model. Two ML models were designed; the results obtained from the proposed model outperformed the existing model found in [43]. The accuracy performance of extreme gradient boost was superb as shown in table 1 with 100% accuracy, precision, recall and f1-score. As for decision tree, its performance is also commendable, 99 percent accuracy and precision while hundred percent for recall and f1-score. From the generated result our proposed models with extreme gradient boost and decision tree algorithms have demonstrated an accurate and reliable performance that is more than what was found in the existing work.

Overall, the results suggest that the XGBoost model is a more effective and reliable model for detecting fake news, and its performance is statistically significantly better than the Decision Tree and KNN models. In general, this research work contributes to the ongoing efforts in addressing the challenges of fake news by utilizing data mining techniques to detect and classify misleading information. The findings have implications for media organizations, social media platforms, and individuals seeking reliable information in the digital age. Limitations of the work emanated from the complex aspect of the fake news detection like Data quality and availability, class imbalance, lack of context, Continuous evolution of fake news, language and cultural barriers, false positives and false negatives, limited domain knowledge and adversarial attacks.

These limitations highlight the challenges and complexities of fake news detection and the need for ongoing research and development to improve the accuracy, effectiveness, and transparency of fake news detection models.

V. CONCLUSION

This work developed an effective approach to identify and categorize fake news articles using data mining techniques. Through the utilization of text preprocessing, feature extraction, and ML algorithms, the proposed approach was capable of distinguishing between real and fake news articles. The performance of the models was evaluated and results demonstrated the effectiveness of data mining techniques in fake news detection and classification. The proposed approach achieves 100% accuracy and performance in distinguishing between real and fake news articles from running the proposed model and obtained results. Media organizations can benefit from incorporating data mining techniques into their fact-checking processes, improving the overall accuracy and reliability of news content. The practical implications of fake news detection involve using specific strategies and tools,

such as fact-checking websites, machine learning algorithms, and social media monitoring, to improve accuracy, efficiency, and credibility, and reduce the spread of misinformation. Fake news research can contribute to restoring public trust in media by developing effective fact-checking methods, improving media literacy, and promoting transparency and accountability in journalism. Also, Social media platforms can utilize these techniques to identify and mitigate the impact of fake news on their platforms, thereby enhancing the trustworthiness of shared information. Moreover, individuals can leverage the outcomes of this study to enhance their media literacy skills and make informed judgments about the credibility of news articles they encounter with. The issue of fake news if not curtailed is causing more harm demanding an imperative action from tech industries and policy makers.

REFERENCES RÉFÉRENCES REFERENCIAS

1. Q. Sheng, J. Cao, X. Zhang, R. Li, Wang, Y. Zhu (2022), "News Environment Perception for Fake News Detection" Key Lab of Intelligent Information Processing of Chinese Academy of Sciences, Institute of Computing Technology, Chinese Academy of Sciences University of Chinese Academy of Sciences.
2. S. Ting (2022), "Automatic Fake News Detection on Twitter" Submitted in Fulfillment of The Requirements for the Degree of Doctor of Philosophy School of Computing Science College of Science and Engineering University of Glasgow.
3. I. Ali, M. Nizam, P. Shivakumara, N. Binti (2022), "Fake News Detection Techniques on Social Media" Wireless Communications and Mobile Computing Volume 2022, Article ID 6072084, 17 page.
4. I. Ahmad, M. Yousaf, S. Yousaf1, M. Ahmad (2020), "Fake News Detection Using Machine Learning Ensemble Methods" Volume 2020, Article ID 8885861, 11 pages.
5. S. Alhabiti a, b, A. Alsalkac, E. Atwell d (2022), "SCUoL at CheckThat! 2022: Fake News Detection Using Transformer-Based Models". The fifth edition of the "Check That! Lab" is one of the 2022 Conference and Labs of the Evaluation Forum (CLEF).
6. E. Johnson1, J. Inyangetoh2, M. Esang3 (2021), "An Experimental Comparison of Classification Tools for Fake News Detection" International Journal of Advanced Research in Computer and Communication Engineering. Vol. 10, Issue 8, August 2021 DOI 10.17148/IJARCCCE.2021.10820 ©IJARCCCE This work is licensed under a Creative Commons Attribution 4.0 International License 135 ISSN (O) 2278-1021, ISSN (P) 2319-5940.

7. U. Sharma, S. Saran, M. Shankar (2020), "Fake News Detection using Machine Learning Algorithms" International Journal of Engineering Research & Technology (IJERT) ISSN: 2278-0181 Published by, www.ijert.org NTASU - 2020 Conference Proceedings.
8. K. Emad, Q. A Al-Haija*, A. Tareef³, M.M. Al-Nabhan⁴ (2022), "Inclusive Study of Fake News Detection for COVID-19 with New Dataset using Supervised Learning Algorithms". International Journal of Advanced Computer Science and Applications, Vol. 13, No. 8, 2022.
9. K. Shu, S. Wang, H. Liu (2017) "Exploiting Tri-Relationship for Fake News Detection" Computer Science and Engineering, Arizona State University, Tempe, 85281, USA.
10. F. Unirio, R. Unirio, A. Unirio (2019) "Can Machines Learn to Detect Fake News? A Survey Focused on Social Media" Proceedings of the 52nd Hawaii International Conference on System Sciences.
11. C. Rodriguez, M. T Keane (2018) "Attention to news and its dissemination on Twitter: A survey" Insight Centre for Data Analytics, School of Computer Science, University College Dublin, Ireland, Computer Science Review 29 (2018) 74–94.
12. A. Bondielli & F. Marcelloni b (2019) "A survey on fake news and rumour detection techniques". a Dipartimento di Ingegneria dell'Informazione, University of Pisa, Largo Lucio Lazzarino, 1, Pisa, Italy. b Dipartimento di Ingegneria dell'Informazione, University of Florence, Italy. m
13. S. Shelke a & V. Attar b (2019) "Source detection of rumor in social network" a Ph.D. Research Scholar, Department of Computer Engineering & IT, College of Engineering, Pune (COEP), 411005, Department of Computer Engineering & IT, College of Engineering, Pune (COEP), 411005, India. Journal homepage.
14. N. Nyow & H. Chua (2019) "Detecting Fake News with Tweets' Properties". Department of Computing and Information Systems Sunway University Selangor.
15. S. Aphiwongsophon & P. Chongstitvatana (2018) "Detection of fake news with machine learning method" 15th International Conference on Electrical Engineering/Electronic, Computer, Telecommunications and Information Technology.
16. V. Rubin, N. Conroy, Y. Chen & S. Cornwell (2016) "Fake News or Truth? Using Satirical Cues to Detect Potentially Misleading News," Proc Second Work Comput. Approaches to Decept. Detect., pp. 7–17, 2016.
17. S. Ray (2017). "common-machine-learning-algorithms"/https://www.analyticsvidhya.com/blog//
18. Economic and Social Research Council. Using Social Mmedia. Available at: https://esrc.ukri.org/research/impact-toolkit/social-media/using-social-media
19. P. Gil (2019), Available at: https://www.lifewire.com/what-exactly-is-twitter-2483331. ASCI-2020 IOP Conf. Series: Materials Science and Engineering 1099 (2021) 012040 IOP Publishing doi:10.1088/1757-899X/1099/1/01204012
20. E. C. Tandoc Jr. et al (2017). "Defining fake news a typology of scholarly definitions". Digital Journalism.
21. K. Sharma, F. Qian, H. Jiang, N. Ruchansky (2019), "Combating Fake News: A Survey on Identification and Mitigation Techniques" University of Southern California MING ZHANG, Peking University YAN LIU, University of Southern California.
22. B. Devi, A. Soni, S. Kapkoti, S, Shankar (2019) "Fake News Detection Based on Machine Learning by using TFIDF" Department of Computer Science and Engineering SRM Institute of Science and Technology, Ramapuram, Chennai, India.
23. Z. Mahid, S. Manickam, S. Karuppayah (2018) "Fake News on Social Media: Brief Review on Detection Techniques" National Advanced IPV6 Centre (Universiti Sains Malaysia) Pulau Pinang, Malaysia.
24. A. Kesarwani, S. Chauhan, A. Nair (2020) "Fake News Detection on Social Media using K-Nearest Neighbor Classifier" School of VLSI and ESD National Institute of Technology Kurukshetra, India.
25. Z. Khanam, B N Alwasel, H. Sirafi¹ and M. Rashid (2021) "Fake News Detection Using Machine Learning Approaches"¹College of Computing and Informatics, Saudi Electronic University, Dammam, KSA² School of Computer Science and Engineering, Lovely Professional University, Jalandhar, India.
26. https://science.sciencemag.org/content/359/6380/1094.summary Science 09 Mar 2018: Vol. 359, Issue 6380, pp. 1094-1096 DOI: 10.1126/science. Aao 2998.
27. Z. Khanam & M.N. Ahsan (2017) "Evaluating the effectiveness of test-driven development: advantages and pitfalls" International. J. Appl. Eng. Res. 12, 7705–7716, 2017.
28. V. Agarwala, H. Sultanaa, S. Malhotraa, A. Sarkarb (2019), "Analysis of Classifiers for Fake News Detection" International Conference on Recent Trends in Advanced Computing 2019, Icrta 2019.
29. J. Chevallier, D. Guégan, S. Goutte (2021). "Is It Possible to Forecast the Price of Bitcoin?" 377–420.
30. K. Poddar., & D, G. B. A. (2019). "Comparison of Various Machine Learning Models for Accurate Detection of Fake News". 1–5.
31. A. Thota, P. Tilak, S. Ahluwalia, N. Lohia "Fake news detection: a deep learning approach", SMU Data Science Review 1 (2018) 10.

32. K. Langin ((2018) "*Fake news spreads faster than true news on twitter*"- thanks to people, not bots, Science magazine.
33. H. Allcott, M. Gentzkow ((2017) "*Social media and fake news*" in the 2016 election, Journal of economic perspectives 31 211–36.
34. J. Gottfried, E. Shearer (2016) "*News use across social media platforms*", [http://www.journalism.org/2016/05/26/news-use-across-social-media-platforms-2016/\(2016](http://www.journalism.org/2016/05/26/news-use-across-social-media-platforms-2016/(2016)
35. C. Silverman, L. Alexander (2016) "*How teens in the balkans are duping trump supporters with fake news*". buzzfeed, 14 november, 2016.
36. DataReportal – global digital insights', Data Reportal –Global Digital Insights. Online. Available: <https://datareportal.com/>. Reached access: 11- Dec-2022.
37. J. M. B (2018), "*Fake News: Real Lies, Affecting Real People*". North Charleston, SC: Createspace Independent Publishing Platform,
38. C. Liu, X. Wu, M. Yu, G. Li, J. Jiang., W. Huang, X. Lu (2019) "*A two-stage model based on BERT for short fake news detection*". Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). 11776 LNAI, pp. 172–183 (2019). https://doi.org/10.1007/978-3-030-29563-9_17_2.
39. X. Zhouk, R. Zafarani (2020) "*A survey of fake news*": fundamental theories, detection methods, and opportunities. ACM Comput. Surv. <https://doi.org/10.1145/3395046>
40. E. K. Qalaja, Q. A. Al-Haija, A. Tareef, A. A. Al-Nabhan. (2022), "*Inclusive Study of Fake News Detection for COVID-19 with New Dataset using Supervised Learning Algorithms*". (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 13, No. 8, 2022.
41. G. Mavridis (2018) "*Fake news and Social Media*". How Greek users identify and curb misinformation online" Media and Communication Studies.
42. Tandoc Jr, E. C. (2019). The facts of fake news: A research review. *Sociology Compass*, 13(9), e12724.