

A MODEL OF CREDIBILITY OF INFORMATION OF SOCIAL PLATFORM**Prajakta Vilas Chiddarwar**

Bharati Vidyapeeth (Deemed to be University), College of Engineering, Pune,
chiddarwarprajakta711@gmail.com

Dr. Sunita Dhotre

Associate Professor, Department of Computer Engineering, Bharati Vidyapeeth (Deemed to be University), College of Engineering, Pune

Abstract:

Today's generation is living in the 21'st century where social media platform is growing rapidly. Increase in the involvement of social media, the number of gateways is getting open to entertaining the user. Many online websites, web applications, and android applications are available over the online platform. Because of these, the credibility of information is introduced on the social media platform. Also fake news, Misinformation, and irrelevant information are critical issues raised on the social platform. Because of these, there is rapid growth in using social platforms like Facebook, Whatsapp, Twitter, Blogs, and many more. Due to this ease of access, people nowadays getting interested in the platform, due to this dissemination of information leads people to spread and search more information.

The dissemination of information may have many ways like fake news, irrelevant data, etc . the impact of this information vary from very low to high, according to the importance of information or search ratio of information. The fake news may contain low-quality data or fake links that spread over a wide or any message to one another. This news can be spread intensively to create an attackable environment to affect the end-user or to damage the status of politicians. Online social media has several fake accounts and that leads to the spread of fake news on the platform. The end-user who is using social media is unaware of the validation and credibility of information which leads to spreading of the fake news from one end to another user.

Keywords: fake news, trust analysis, the credibility of information, machine learning

Introduction:

Using an online social platform, users can easily connect to millions of people and can communicate over the network. People can upload or download content from social media. People use social media for gathering relevant information and also spreading information or news over the media. In recent years, the spreading of fake news and Misinformation is growing rapidly. This information can be regarded by politicians, personal, actors, sportspersons, and any high demanded person. Therefore the credibility of information is a central issue in the recent era of social media.

The major impact of fake news was gathered in 2016, in US Presidential Election. At that time number of fake news and irrelevant information were spread over the network which leads to a major impact on the election. After analyzing the various news, the impact factor of

fake news was high than the relevant information which leads to getting diverse decisions in an election. This result shows that the ease of use of the social platform is getting in the wrong manner to take advantage of getting information.

Uploading and downloading information from social media is very easy as well as very cheap. When an end-user sends a message on the web application, the receiver will get the message fast but at the same time the receiver is unaware of the trustworthiness of the message and because of this the dissemination of information is formed on social media.

Detecting fake news on social platforms is now a day new task as the spread rate of fake news is very high. If a sender sends a message about any politician or influencer, the information can spread from one end to another end widely and can create a big issue on social media, detection of fake news concept was introduced to measure the credibility of information that spread over the internet.

1.1 Credibility of Information:

Credibility is the term of correctness, how much the content of information is correct and how it is trustworthy. Credibility also refers to the quality and accuracy of the message or context that is shared on a social platform. Credibility proposed that the message that is shared on a platform has to be accurate and un-bias as well as there should not be any manipulated data shared on the social platform. The source and destination should be constant in nature. The credibility of information is affected in three terms, the source of the message, the message itself, and the platforms on which the message is shared.

Before online social platforms, credibility was checked in an offline context However in the online platform era, credibility is checked in an online context. Because of the ease of sharing information, the dissemination of misinformation is widely increased and the context has to be checked for correctness. Also, online context can be manipulated and can be easily altered which makes the context vulnerable. The credibility help to restore the accuracy and quality of the context that is shared on the social platform.

Literature Survey:

“In this section, we briefly review the related work on fake news detection system and their different techniques.

In this paper [1] , The results of a fake news identification study that documented the work of a fake news classifier are shown. Textblob, Natural Language, and SciPy toolkits were used to develop a new fake news detector. It is an advantage to have. It uses natural language processing. Second, fake news detection based on attribute classification. Cons: This is a arduous process.

This treatise [2] introduces a dataset that contains both fake news and real news, and conducts various experiments to sort out fake news detectors. The downside is that it uses a limited dataset.

This paper [3] proposed a distributed framework to implement the proposed truth discovery scheme using Work Queue in an HTCondorsystem. Advantages is it Finds trustworthy

information on Social media and second .Proposed truth discovery scheme using Work Queue in an HTCondor system and disadvantages is Accuracy is low.

Paper [4] Studied various detection techniques i.e. content based, social context based and hybrid based. Advantages is Proposed content-based, social context-based and hybrid-based methods and disadvantages is only survey state of the methods.

This paper [5] Present a new fake news detection model using unified key sentence information which can efficiently perform sentence matching between question and article by using key sentence retrieval based on bilateral multi perspective matching model. Advantages is Implement natural language processing using key sentence retrieval and disadvantages is Fake news detection accuracy is low.

This Paper [6] classifies fake news messages from Twitter posts using hybrid of convolutional neural networks and long-short term recurrent neural network models. The advantage is the implementation of hybrid CNN and RNN models, which are much more accurate. The downside is just the headline of the tweet.

This paper [7] Compares 2016 news with other sources. The advantage is to first detect the 2016 fraudulent elections that spread through social media. This treatise provides the theoretical and empirical background for shaping this argument. The disadvantage is. The dataset used is limited and is limited to 2016 news only.

This paper [8] presents a new approach for detecting fake news using a naive Bayes classifier. It uses an implementation of the Naive Bayes machine learning algorithm, but with less accuracy.

This paper [9] introduced the basic concepts and principles of fake news in both traditional and social media. The detection phase reviews existing approaches to detecting fake news from a data mining perspective, such as feature extraction and model building. The advantage of this paper is that it investigated the issue of fake news by reviewing existing literature in two phases. H. characterization and detection. However, it uses static data.

This study [10] contributes to scientific knowledge of the impact of different types of media use interactions on political impact. The advantage is that multiple news sources are used to detect fake news, and the disadvantage is that it focuses only on political data.

Proposed Work:

In the new era of online social platforms, the usage of an online web application is growing rapidly, and the spreading of information is increased from one end to another. Therefore, the trustworthiness of information is raised on social media. The credibility of information concept is proposed to detect fake news from the online social platform. As the news that is received on one end is collected from different resources and the data correctness and quality of data are also manipulated. Online social media such as Facebook, Twitter, WhatsApp, Blogs, etc. are also doubted the accuracy of information. To solve this problem, we have special annotators who are experts in the fields of "NLP (Natural Language Processing)" and "Text Extraction". It extracts the context we get and retrieves the news module so that we can compare it to real news.

In this paper, flow of model is as shown in below:

Source: Author or publisher of the news article.

Headline: Short title text that aims to catch the attention of readers and describes the main topic of the article

Body Text: Main text that elaborates the details of the news story; there is usually a major claim that is specifically highlighted and that shapes the angle of the publisher

System Diagram:

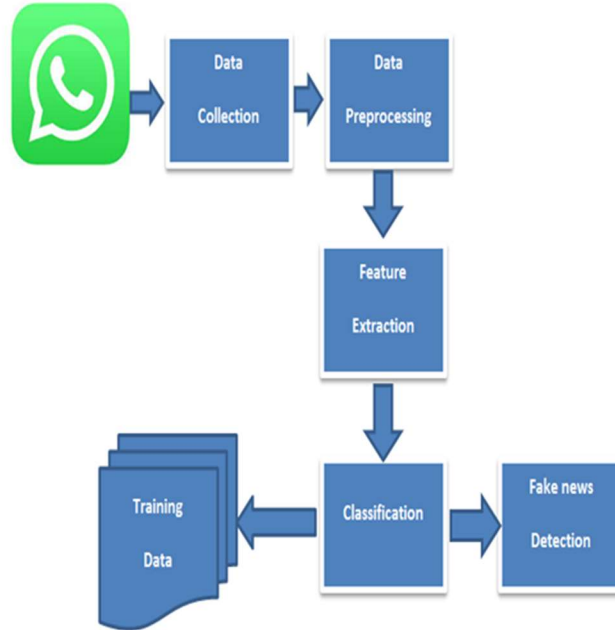


Fig 1. System Work Flow

System Work Flow:

Data Collection:

One of the hardest problems to solve in machine learning has nothing to do with complex calculations: it's a matter of putting the right data sets in the right organizations. Obtain accurate information regarding aggregating or discriminating information related to predicted outcomes. Selecting the right dataset for machine learning is crucial to making the AI model work with the right approach, while choosing the right quality and quantity of data.

As this project is a real-time web application, In this project I have used user define data set to get the accuracy of the message and it is handled by the admin of the model.

Data Pre-processing:

In the machine learning process, preprocessing data is the step of transforming or encoding the data so that it can be easily analyzed by the machine. In other words, the characteristics of the

data can now be easily interpreted by algorithms. For this fake message detection, the most important thing you need to do is preprocessing. First, records are collected from a variety of sources, so you need to remove unnecessary information, convert it to lowercase, and remove punctuation marks, symbols, and stop words.

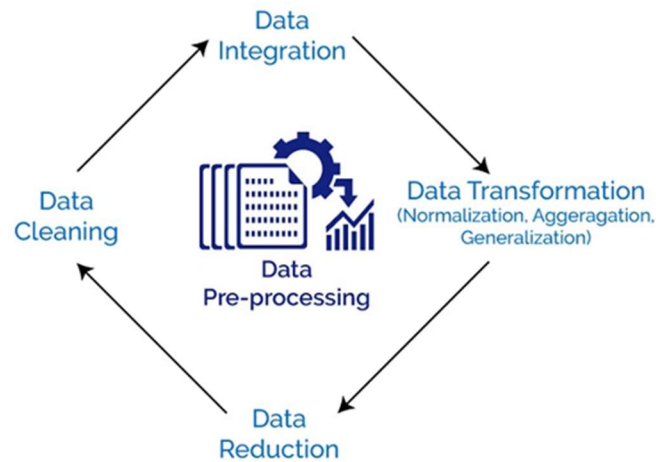


Figure 2: Data Preprocessing

Feature Extraction:

Feature extraction improves the accuracy of the trained model by extracting features from the input data. In this phase of a typical framework, you reduce the dimensions of your data by removing redundant data. Of course, it speeds up training and reasoning. The feature extraction method retains the newly created features by performing a combination and transformation of the original feature set.

Training and Testing:

To create a valuable training set, you need to understand the problem it is doing. For example, what machine learning calculations do and what yields are expected. Machine learning works regularly with two sets of information: training and testing. Each needs to randomly test a wider range of information. The main set used is the largest of the two training sets. When you run a training set through a machine learning system, the network shows you how to weight different highlights and change the coefficients, depending on the potential to limit the resulting errors. These coefficients, also called parameters, are contained in the tensor and are collectively referred to as the model because they encode the model of the information to be trained. These are important insights from the preparation of machine learning systems. The second set is the test set. It acts as a sign of approval and is only used at the end. Once you have prepared and set up the information, you can test the neural network against any of this last test. The results produced must ensure that the network recognizes the image correctly or at least remembers the image at the [x] level. If the exact predictions are not met, return to the training set and check for any errors that have occurred. It is okay to include the correct records and the system will run smoothly.

Algorithm Used for Classification:

This section is about classifier training. Various classifiers have been explored to predict classes of texts, and in particular, four machine learning algorithms have been studied: – Multinomial Naïve Bayes Passive Aggressive Classifier and Logistic regression An implementation of these classifiers were created using the SciKit Learn Python library.

Introduction to the Algorithms:**Naïve Bayes Classifier:**

Naive Bayes is a supervised learning algorithm based on Bayes' theory and used to solve classification problems. It is primarily used for text classification containing the multivariate training data set. The Naïve Bayes Classifier is one of the simplest and most efficient classification algorithms to help you build fast machine learning models that can make fast predictions. A probabilistic classifier. In other words, it makes predictions based on the probability of an object.

Random Forest Algorithm:

Random Forest is the brand name for Decision Tree Ensemble. A random forest has a set of decision trees (called "Forests"). To classify a new object based on its properties, each tree provides a classification and the tree says to "vote" for that class. Forest chooses the classification with the most votes (out of all trees in the forest). Random Forest is a classification algorithm consisting of decision trees. It uses bagging and feature randomness when building each individual tree to try to create an uncorrelated forest of trees whose prediction by committee is more accurate than that of any individual tree. Random forest, like its name implies, consists of a large number of individual decision trees that operate as an ensemble. Each individual tree in the random forest spits out a class prediction and the class with the most votes becomes our model's prediction. The reason that the random forest model works so well is: A large number of relatively uncorrelated models (trees) operating as a committee will outperform any of the individual constituent models. So how does random forest ensure that the behaviour of each individual tree is not too correlated with the behaviour of any of the other trees in the model? It uses the following two methods:

- Bagging (Bootstrap Aggregation) - Decision trees are very sensitive to the trained data. - Small changes to the training set can significantly change the tree structure. Random forests take advantage of this by randomly selecting each individual tree from a data set through replacement, allowing the creation of another tree. This process is called bagging or bootstrapping.
- Feature Randomness— When splitting nodes in an existing decision tree, consider all possible features and choose the one that provides the greatest separation between observations from the left and right nodes. In contrast, each tree in any forest can only be selected from a random subset of features. This results in more variability between the trees in the model, which ultimately leads to lower correlations and diversification between the trees.

Logistic Regression:

This is a classification, not a regression algorithm. Used to evaluate discrete values (binary values such as 0/1, yes/no, true/false) based on a given set of explanatory variables. Simply put, fit the data to a logit function to predict the probability of an event occurring. Therefore, it is also called logit regression. Because we are predicting probabilities, the output values are between 0 and 1 (as expected).

Passive Aggressive Classifier:

The passive-aggressive algorithm is the online algorithm. Ideal for classifying large data streams such as Twitter. Easy to implement and very fast. For example, it works by learning and then discarding. These algorithms remain passive on correct classification results, while becomes aggressive in case of miscalculation, updating and correcting. Unlike most other algorithms, it does not converge. Its purpose is to make an update that corrects for the loss by making very small changes to the proportions of the weight vectors.

Random forest algorithm

The Mathematics behind Random Forest

Gini Index

Method to split out the data is the Gini index; it checks the impurity or purity of data it is used in the CART (Classification and Regression Tree) algorithm like Decision Tree.

It generates a binary split, which is then used by the CART algorithm. An attribute is low Gini index is preferred as the root node

Formula to calculate Gini index is:-

$$\text{Gini Index} = 1 - \sum_j P_j^2$$

Information Gain

Information gain is calculated using the attributes entropy and entropy in the data set, which shows us how much information a feature offers us with a class.

$$\text{Entropy} = - \sum p(x) \log p(x)$$

The amount of impurity or randomness in the data is measured by entropy. It is used to determine the decision tree's root node for data splitting.

Formula to calculate information gain

$$\text{Information Gain} = \text{Entropy}(S) - [(\text{Weighted Avg}) * \text{Entropy (each feature)}]$$

Highest the information gain we select that feature as the root node

Regression Problems

RFA to solve regression problems that time you are using the mean squared error (MSE) value to how your data branches from each node.

$$\text{MSE} \frac{1}{p} \sum_{i=1}^p (Q_i - R_i)^2$$

Where,

P= number of data points,

Qi= Value returned by the model

Ri= Actual value of data point i

Pseudocode of Random Forest Algorithm

Precondition: Training Set : S

Function Random forest(S,D)

for i=1 to S do

Sample data randomly replace with Di

Create root R contains data of Di

Build Tree (Ri)

end for

Build Tree (R):

If R have instance of single class then

Return

else

Select randomly possible splitting feature in R

Select feature f with highest information

Create n child node of R where F has n possible values (n1,n2,...nn)

For i=1 to n do

Set the content of Ri to Di where Di is all instance in R match with Fi

Call Build Tree (Ri)

end for

end if

V.RESULTS AND DISCUSSION**Step I : Home Page**

See Figure 5.1 for our system's home page, which contains our false news detection system's home page.

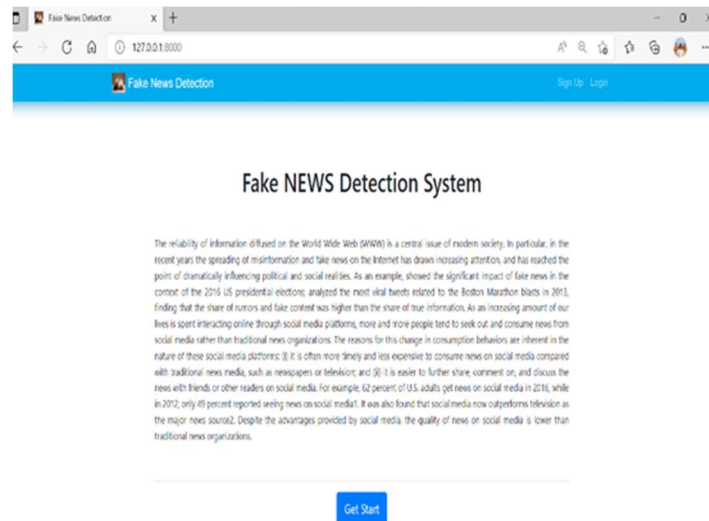


Figure 3: Home Page

Step II: Signup Page

As indicated in Figure 5.2, the parameters username, email, password, and confirm password will be displayed, and the user should sign in to our system using this information.

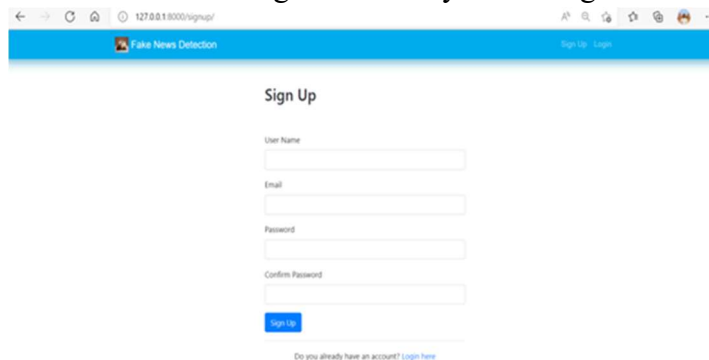


Figure 4: Signup Page

Step III: Login Page

Our third step is to take the user to a login page where they may enter their username and password to access our system.

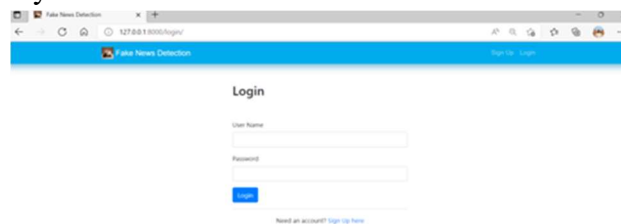


Figure 5: Login Page

Step IV: Input Post

In the fourth stage, the user/we post the news that Mahindra Singh Dhoni has announced his retirement from all forms of cricket, as seen in Figure 5.4.

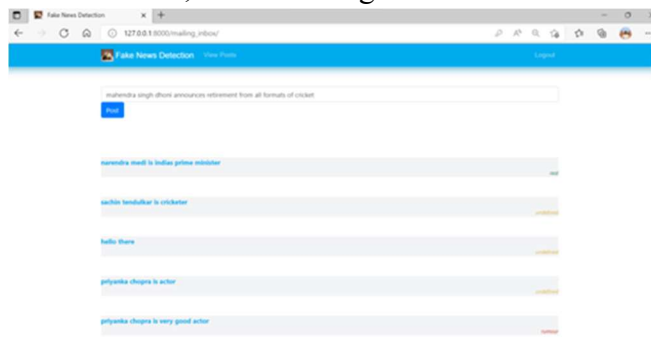


Figure 6 : Input Post

Step V: Undefined Post

Because we don't know if the post is real or fake in Undefined Post, it will appear as undefined post in the figure below as soon as we enter fresh data into the system. Mahindra Singh Dhoni has announced his retirement from all forms of cricket, and the news will be marked with an undefined tag.

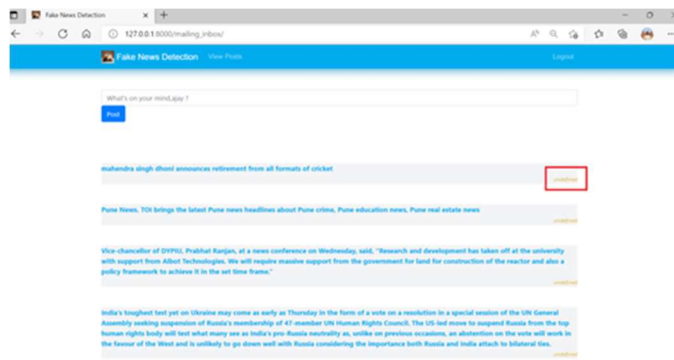


Figure 7: Undefined Post

Step VI: Admin to select post is real or rumour

We can see the admin panel, admin panel is used to select post is real or rumore. Using admin panel admin decide the news is real or not and after deciding admin will submit the answer that is about new is real or not.

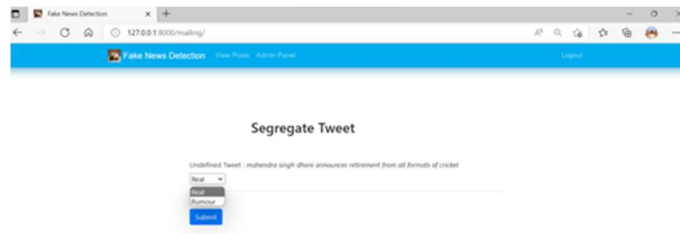


Figure 8: Admin to select post is real or rumour

Step VII: Similar Post get atomically call segregated as real

This is last step of our system in this step system will tell about news is real or not and also segregate similar post according to current post and decide news is real or not. As seen in figure 5.7 it display the various news and according to our system it will also display the news is real or not.

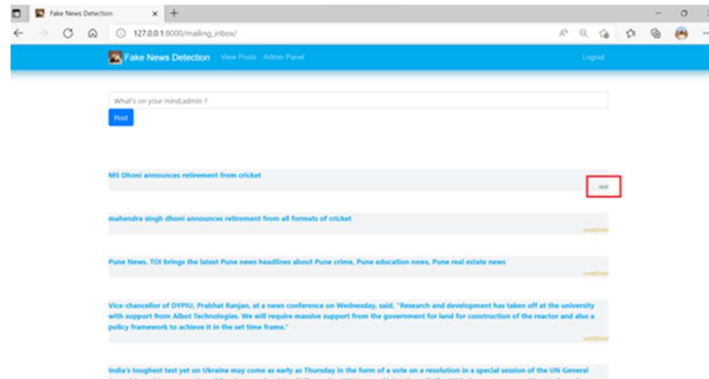


Figure 9: Similar Post get atomically call segregated as real

Confusion Matrix

	Class 1	Class 2
Class 1	1530	4
Class 2	5	1535
Total for Class	1535	1539

Figure 10 : Confusion Matrix

The confusion matrix Class 1, Class 2 training modules can be seen in the diagram above. In Class 1, the input photos are 1535, and we achieved accuracy of 99.71 % and precision of 1.0 % while training the classifier as a train with the supplied input database. Because the 1535 classifier failed to classify 4 photos as an output form of a class 1, recall was reduced to 1.0%, and F1 score was also reduced to 1.0%.

In Class 2, the input photos are 1540, and we achieved accuracy of 99.71 % and precision of 1.0 % while training the classifier as a train with the given input database. As a result of the

1539 classifier failing to detect 5 photos as an output form of a Class 2, recall has been reduced to 1.0 %, and F1 score has been reduced to 0.99%.

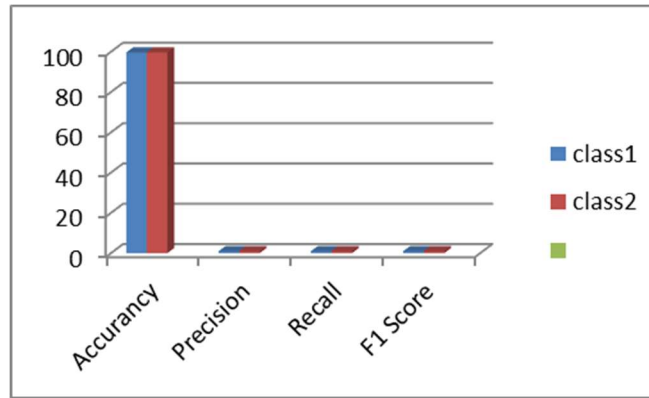


Figure 11 : Performance analysis graph

The system gives 99.71 % accuracy using RFA.

Comparison of Accuracy, Precision, Recall and F1 Score for all four classifier.

Classifier	Accuracy	Precision	Recall	F1 Score
Naïve Bayes	87.62	0.81	0.9	0.85
Logistic Regression	73.5	0.64	0.72	0.68
Random Forest	99.71	1	1	1
Decision Tree	89.11	0.89	0.88	0.89

Table 1: Comparison Table

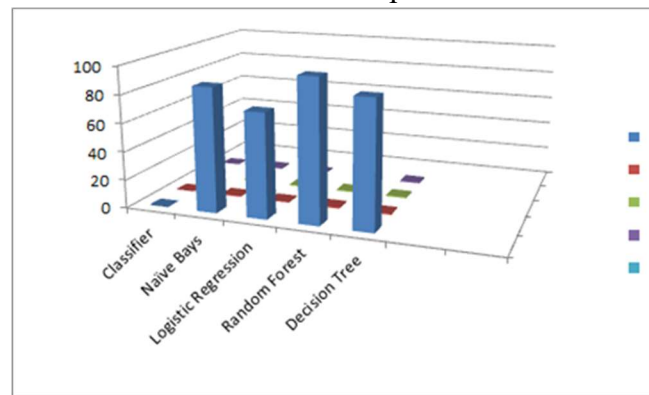


Figure 13: Graph Representation for all Four Classifier

Shown in figure 13 it shows that RFA gives more accuracy than any of three.

VI.CONCLUSION

People are increasingly turning to web-based media for news rather than traditional media as web-based media becomes more popular. However, the internet has been used to spread misleading information, causing harm to both individual customers and society as a whole. News by examining existing writings in two stages: representation and recognition. Throughout the representation phase, we explain the basic ideas and criteria of FNs in traditional and web-based media. Throughout the localization phase, we look at current FNs detection from an

information mining perspective, including highlight extraction and model structure. Datasets, evaluation measures and possible future paths in simulated identification research are also discussed, as is expanding the discipline to cover a broader range of applications. The system gives 99.71 % accuracy which is better than any other algorithm.

REFERENCES

- Fazlullah Khan, Abdul Wali Khan University Mardan, Pakistan Syed Tauhidullah Shah, University of Calgary, Canada Nabeela Awan, Nagaoka University of Technology, Japan,"Special Issue on Advanced Aspects of Machine Learning Algorithms for Scientific Programming:"2022
- Ridhima Mehta,"Applying fuzzy logic for multicriteria performance analysis of social media networking" 2022
- Vishal jain, jyotir moy chatterjee, ankita bansal, utku kose and abha jain, volume 13 in the series de gruyter frontiers in computational intelligence "computational intelligence in software modeling"2022
- Terry Traylor, Jeremy Straub, Gurmeet, Nicholas Snell" Classifying Fake News Articles Using NLP to Identify In-Article Attribution as a Supervised Learning Estimator"2019.
- Rohit Kumar Kaliyar" Fake News Detection Using A Deep Neural Network"2018.
- Daniel (Yue) Zhang, Dong Wang, Nathan Vance, Yang Zhang, and Steven Mike" On Scalable and Robust Truth Discovery in Big Data Social Media Sensing Applications" 2018.
- ZaitullradahMahid,SelvakumarManickam,ShankarKaruppayah" Fake News on Social Media: Brief Review on Detection Techniques" 2018.
- Namwon Kim, DeokjinSeo, Chang-Sung Jeong" FAMOUS: Fake News Detection Model based on Unified Key Sentence Information" 2018.
- Oluwaseun Ajao, DeepayanBhowmik, ShahrzadsZargari" Fake News Identification on Twitter with Hybrid CNN and RNN Models"2018.
- Hunt Allcott Matthew Gentzkow" SOCIAL MEDIA AND FAKE NEWS IN THE 2016 ELECTION" 2017
- MykhailoGranik, VolodymyrMesyura" Fake News Detection Using Naive Bayes Classifier"2017
- Kai Shu , Amy Sliva , Suhang Wang , Jiliang Tang , and Huan Liu" Fake News Detection on Social Media: A Data Mining Perspective"2016
- MeitalBalmas" When Fake News Becomes Real: Combined Exposure to Multiple News Sources and Political Attitudes of Inefficacy, Alienation, and Cynicism" 2014