

Fake News Detection using Data Science Approach

**Prof. Mrs. Rupali Shinganjude, Ayush Khadgi, Harsh Watkar,
Harshad Bagde, Brajesh Dayare, Dipak Choudhary**
Department of Information Technology
Priyadarshini Bhagwati College of Engineering Nagpur, India

***Abstract:** Information plays an vital role in an every aspects. Thus the misleading in information is known as fake news. This news create the controversy among the people and tragic events may take place. Fake news has been problem ever since the internet boomed. The very network that allows us to know what is happening globally is perfect breeding ground for malicious and fake news. Combatting this fake news is important because the world's view is shaped by information. This research paper proposed a reliable system which is able to classify the fake news using the NumPy, pandas Sklearn, Vectorization in python, random forest classifier, decision tree classifier logistic regression classifier, gradient boosting classifier, regular expression library in python String library in python, seaborn package.*

Keywords: Fake News

I. INTRODUCTION

In today's era where the internet is one of the major source of obtaining information regarding the current or previous scenarios.[8]Modern life has been quite convenient and the people of the world have to thank the immense contribution of the internet technology for the communication and information sharing. Fake news spread mostly through social media. Fake is threat to the politics, finance, education, democracy, business. Although fake news is a previously ongoing problem but nowadays human believe more in social media which leads to believe in fake news and then spread of the same fake news.

Determining the authenticity of content is critical to preventing the spread of fake news, and it is of paramount importance the we confront and try to eliminate the problem of fake news. The approaches we were used with the help of python libraries are String , RE(regular expression library in python). This library extracts textual part from the article which is going to provide by users. Regular expression is a pattern that the regular expression engine attempts to match in input text. A pattern consists of one or more characters literals, operators, or constructs.

II. LITERATURE REVIEW

Our approach is to build a system that is capable of deciding real and fake news. Jiawei Zhang, Bowen Dong [1] the problem you are going to face while handling dataset of news or creating it, before this formal definition and formulation of the problem is required necessary before the studying the problem. To resolve these challenges, the paper introduces a new fake news detection framework, the name is given as FakeDetector. Syed Manzoor, Dr.Jimmy Singla [2] the information about the novel method and tool for detecting fake news. It uses the text preprocessing, encoding of the text, extraction of the characteristics, support vector machine. Steaming and analysis the text by removing stop words and special characters. Using sack of words and n-gram the TF-IDF. The source of a news, its author, the date and the feeling given by the texts as feature of a news. A supervised machine learning algorithm that allows the classification of new information. Z Khanam , B N Alwasel , H Sirafi and M Rashid [3] propose to create the model using different classification algorithm. The model will examine the unseen data, the result will be plotted. The model that detects and classifies fake articles and can be used and integrated with any system for future use. This paper includes various machine algorithm like linear regression, random forest, XGboost, Naïve bayes, K-nearest neighbor(K-NN), decision tree, support vector machine(SVM). Text data requires preprocessing before applying classifiers on it, we will clean noise, using NLP(natural language processing) for POS(part of speech) processing and tokenization of words. Vidhi Singrodia, Anirban Mitra[4] It is a procedure which automates the web data extraction instead of manually copying it. The data in PDF documents is in an unstructured format whereas the HTML document

has the data in a structured format. consequently, there is a need for automation of the procedure, apart from the tedious monotonous and faulty duties of physical data retrieval.

Random forest classifier

The goal of ensemble method is to combine the prediction of several base estimators built with a given learning algorithm in order to improve generalizability / robustness over single estimators. [11]Forest classifier have to be fitted with two arrays: a sparse or dense array X of shape($n_samples, n_features$) holding the training samples and array Y of shape the target values (class labels). For classification tasks, the output of the random forest is the class selected by most trees. However, they are seldom accurate. Random forest are a way of averaging multiple deep decision trees, trained on different parts of the same training set with goals of reducing the variables.

[12]Samples with replacement, n training examples. From X, Y ; call these X_b, Y_b . train a classification or regression tree F_b on X_b, Y_b . after training, predictions for unseen samples X' can be made by averaging the predictions from all the individual regression trees on X' .

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(X')$$

Or by taking the majority vote in the case of classification trees. The bootstrapping procedure leads to better performance because it decreases the variance of the model, without increasing the bias.

Decision Tree classifier

Classification is a two-step process, learning step and prediction step, in machine learning. [13]In the learning step the model is based on given training data. In the prediction step, the model is used to predict the response for given data. Decision tree is one of the easiest and popular classifications algorithm to understand and interpret. decision tree algorithm is the member of the supervised learning algorithms family. This algorithm also solves the regression and classification problem too. the objective behind using this algorithm is to create a training model that can use to predict the class or value of the target variable by learning simple decision rules inferred from prior data (training data). It is an non-parametric supervised learning algorithm. Data comes in the record form $(x, Y) = (x_1, x_2, x_3, \dots, x_n, Y)$.

$$H(x) = - \sum_{i=1}^n p(x_i) \log(p(x_i))$$

Logistic regression classifier

Logistic regression is a statistical method used to analyze and model the relationship between a binary (yes/no) dependent variable and one or more independent variables. [14]It is a type of regression analysis that is used to model the probability of a certain outcome or event occurring based on one or more predictor variables. The output of a logistic regression model is a predicted probability of the dependent variable taking the value of 1 (i.e., the event occurring), given the values of the independent variables. This probability is estimated using a logistic function, which maps any real-valued input to a value between 0 and 1. The logistic regression model is based on the assumption that the relationship between the dependent variable and the independent variables is linear on the logit scale. The logit transformation is applied to the predicted probability of the event occurring, which produces a linear relationship between the independent variables and the transformed outcome. [15]The logistic regression model then estimates the coefficients of this linear relationship using maximum likelihood estimation. Logistic regression is commonly used in various fields, including medical research, social sciences, marketing, and finance, among others. It is particularly useful when the dependent variable is binary and the independent variables are continuous, categorical, or a combination of both. The logistic regression model can be expressed mathematically as:

$$\text{logit}(p) = \ln(p/(1-p))$$

$= \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ where p is the probability of the event occurring, \ln is the natural logarithm, β_0 is the intercept term, β_1, \dots, β_k are the coefficients for the independent variables x_1, \dots, x_k , respectively. The logit function is the link function that transforms the predicted probability of the event occurring to a linear combination of the independent variables. The logistic regression model estimates the coefficients $(\beta_0, \beta_1, \dots, \beta_k)$ using maximum likelihood estimation, which involves finding the values of the coefficients that maximize the likelihood of the observed

data given the model. Once the coefficients are estimated, the logistic regression model can be used to predict the probability of the event occurring for a new observation by plugging in the values of the independent variables into the model and solving for p . In binary classification, precision is a metric that measures the proportion of true positive predictions out of all the positive predictions made by the model. It can be defined mathematically as:

$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$\text{Recall} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

The F1-score can be calculated using the precision and recall values obtained from the confusion matrix of the logistic regression model. For example, if the confusion matrix is:

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

Then the precision and recall can be calculated as:

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

And the F1-score can be calculated as:

$$\text{F1-score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

III. MODULES DISTRIBUTIONS

This section gives the instructions about the methods, algorithms, libraries used in the overall project or a system. Our system consist four modules i.e., first and second module works on the textual extraction which means this modules makes the data in more readable format. The regular expression libraries are working in these modules. [9]The third module has train_test_split method which trains and tests the readable format data. the last module of our system collaborate with our four classifiers random forest classifier, decision tree classifier, logistic regression classifier, gradient boosting classifier.

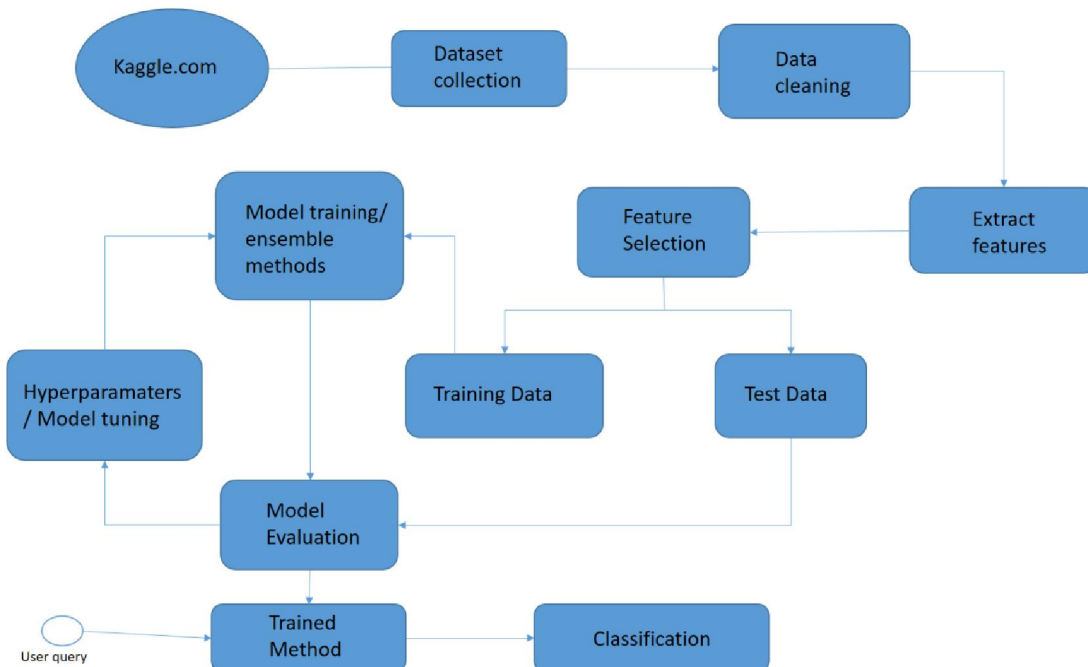


Fig: flowchart diagram for fake news detection.

First thing is going to happen that previous dataset from kaggle. We have to clean the data using the regular expression library in predefined for the cleaning task. So, removing all the special characters and unwanted space, quote, and many more. After this, we get the purely the textual data that we want for the training aur testing purpose. Now we are going to separate the text attribute from the dataset. For this we drop other attribute present in the dataset. So we have only the text part than we trained the text part and label them as 1 and 0 which is binary. The label 1 represent True and the label

0 represent False. Now these label dataset going to merge and form the only dataset with both the label of True and False. Dataset goin to divide in two one is for training and other is for testing. testing part is going for the model evaluation. Model evaluation further process this testing dataset to hyperparameters/model tuning. This process further move the data to model training and ensemble method and after this it returns to the model evaluation and proceed to trained module section and lastly it goes to the classifier also for classification. in between trained methods users request is also manages to give the output.

IV. CONCLUSION

The vast majority of tasks are done through the internet in the twenty-first century. Applications like Facebook, Twitter, and online news stories are replacing newspapers, which were formerly preferred as printed copies. The forwards on WhatsApp are another important source. Fake news, a problem that is only becoming worse, complicates matters and seeks to influence or sway people's attitudes toward using digital technologies. When a person is confused by the genuine news, one of

two things may occur: First, they may begin to believe that their observations of a certain subject are accurate. Therefore, in order to stop the issue, we created a system called Fake News Detection that classifies user input as either real or false. Various NLP and machine learning techniques must be employed to do this. A suitable dataset is used to train the model, and its performance is also tested using a variety of performance measures. The classification of news headlines or articles is done using the best model, or the model with the highest accuracy. Our best model turned out to be Logistic Regression, which had an accuracy of 65%, as is shown above for static search. In order to improve the performance of logistic regression, we therefore employed grid search parameter optimization, which provided us with an accuracy of 75%. Therefore, we can state that there are 75% possibilities that a user would successfully classify a given news story or its headline according to its true nature if they submit it to our model. Users can research the effectiveness of news articles, keywords, and websites online. Cycle after cycle, the accuracy of the dynamic system improves to 93%. I would like to create my own dataset that is regularly updated with the latest information. We use web crawlers and online databases to store the latest information and live news.

REFERENCES

- [1]. Jiawei Zhang, Bowen Dong , Philip S. Yu IFM Lab, Department of Computer Science, Florida State University, FL, USA ArXiv:1805.08751v2 [cs.SI] 10 Aug 2019
- [2]. Syed Manzoor, Dr. Jimmy Singla, Nikita Proceedings of the Third International Conference on trends in electronic and Informatics (ICOEI 2019)
- [3]. Z Khanam , B N Alwasel , H Sirafi and M Rashid IOP Conference Series: Materials Science and Engineering Z Khanam et al 2021 IOP Conf. Ser.: Mater. Sci. Eng. 1099 012040
- [4]. Akshay Jain, Amey Kasbe. 2018 IEEE International Students' Conference on Electrical, Electronics and Computer Sciences.
- [5]. Vidhi Singrodia, Anirban Mitra 2019 International Conference on Computer Communication and Informatics (ICCCI -2019), Jan. 23 – 25, 2019, Coimbatore, INDIA
- [6]. Daryna dementieva, alexander panchenko. Fake news detection on multilingual .2020 IEEE DOI june 22, 2021
- [7]. Fuad mire Hassan, mark lee political fake news detection via multistage feature-assisted neural modeling. Dec 21, 2020.
- [8]. Shikun Iyu, dan chia-tien lo. Fake news detection by decision tree. 2020 IEEE. June 02, 2021.
- [9]. Nihel Fatima baarir, abdelhammid djeflal, fake news detection using machine learning approaches, 978-1-6654-4084-4/21, 2020 IEEE, may 24, 2021.
- [10]. Krishna A N, member IEEE, identification of fake news using machine learning, 978-1-7281-6828-9, sep 17, 2020.
- [11]. Random forest classifier <https://www.analyticsvidhya.com/blog/2021/06/understanding-random-forest>
- [12]. Random forest <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>.
- [13]. Decision tree classifier https://en.wikipedia.org/wiki/Decision_tree_learning.
- [14]. Logistic regression <https://www.ibm.com/in-en/topics/logistic-regression#:~:text=Resources,What%20is%20logistic%20regression%3F,given%20dataset%20of%20independent%20variables>.

[15]. Logistic regression mathematical statement <https://www.analyticsvidhya.com/blog/2015/11/beginners-guide-on-logistic-regression-in-r/>