

Fake News Detection using ML

Dr. Jayakumar V^a, Niket kumar^b, Navneet Himanshu^c

^aProfessor School of Computing Science & Engineering (SCSE), Galgotias University

^{b,c} School of Computing Science & Engineering (SCSE), Galgotias University

^ajayakumar.v@galgotiasuniversity.edu.in, ^bniketkumar838@gmail.com, ^csonugupta2196@gmail.com

Article History: Received: 11 January 2021; Revised: 12 February 2021; Accepted: 27 March 2021; Published online: 10 May 2021

Abstract: This Project “Fake News Detection” works on the applications of Natural Language Processing (NLP) techniques that recognizes the 'fake news', that is deceptive news stories which comes from the unidentified sources. During this systematic review, the factors that results in the spreading of fake news and information have been provided. In this report, the identification of the basic cause which results in the spreading of fake news are performed which may result in the break of fake news among public domain. In order to conquer the social media platform from the rapid spread of fake news, firstly we should know the reason and intention behind the spreading of fake news. Therefore, this review takes associate in early initiative to find the major reason which lead to the expansion of fake news among public domain. The main aim of this review is to find out with what intention and why people unknowingly share information which may be false and to presumably facilitate in detection of fake news before it spreads. There the model should be build which support a count vectorizer or a (TFIDF) Term Frequency Inverse Document Frequency matrix, will solely get you up to now. However sometimes these following models did not consider the important qualities like ordering of word and context. It may be possible that 2 articles whose word counts may be similar are totally alter in their meanings

Keywords: Fake News, Misinformation, Social Media, Classification

1. Introduction

The emergence in the swift increase of internet users and also the fast acquisition of social media platforms such as twitter and facebook sealed the method for circulation of information that has never been witnessed within the human history earlier. Fake news refers to false or misleading information content whose source cannot be verified. Besides different instances, news retailers gained from the widespread use of different media network platforms by delivering updated news to their subscribers through apps, and different digital platforms they can be either facebook or website or whatsapp or twitter, blogs, social media feeds, and other digital media platforms. Social media platforms such as twitter, facebook, instagram, whatsapp, etc are considered so much influential when it comes to news feed. This became quiet for customers to accumulate the most recent news very quickly. These social network platforms are gaining such huge popularity because in today's time they provide an edge to their users that they can express their feelings as well as discuss together and represent their thought in front of society, they can share their opinion over the topics like health, poverty, education and literacy. But these social media platforms such as facebook, whatsapp, twitter can also be used in negative manner where they can be used to spread fake news across the society, negativity among the youth which can lead to riots as we have seen many cases and these cases are increasing every year exponentially. These fake news on one place can cause mass destruction but for somebody present in between of us can earn money by selling such news.

Fake news is not only limited to small point they got so much popularity during the time of USA election 2016 as well as 2021. Huge mass of population was sharing information among each other through social media mainly facebook without knowing the facts which lead to serious issues. These fake news articles were not just covering politics but much more. They have alternative domains too such as health, sports, science as well as lifestyle too. Financial market is the worst affected area from fake news and articles, Over here a small fake news can be disastrous which can lead to a halt on the market which will be loss making. This world is formed by the data or information which we can digest. There is a proof that customers act angrily over the news which later found out to be fake or false or incorrect. Recent fake news which was spreading was related to the topic Novel corona virus. There were so many misinformation regarding the behaviour of virus, its spread as well as many aspects which were fake but provided to very large chunk of population.

But we are bit lucky that there are various techniques and procedures in computer science which we can use to mark those articles which are present completely textual content. There are websites such as “PolitiFact” which can be used for fact check of information or article. Researchers have maintained depositories which contains the inventory of websites that were recognised as the source of fake news or ambiguous and fake. These techniques or methods are helpful in training of many machine learning algorithms in a more efficient and effective manner. But, the matter of concern with these methods or resources is that our or human experience is needed to identify that whether the given articles or websites is spreading misleading information or fake. There are plenty of fake news busters or fact checking websites which consist of articles from specific domains, such as politics and doesn't seem to be generalized to spot fake news articles from different domains like sports, and technology.

Many researches have mainly targeted upon the detection and classification of fake news on different social media platforms such as whatsapp, instagram, Twitter and Facebook. At abstract level or top level, fake news is determined into completely different types; the knowledge is then dilated to generalize (ML) machine learning models for multiple domains .

1.1 Contribution

In the compilation of fake news, we need to consider multiple instances where both supervised and unsupervised machine learning is used for the classification of texts. Though mostly the literature emphasizes upon certain specific domains like political domains. On getting the result we can see our algorithm works better for a particular type of domain and gives optimum results but when it focuses on different type of domains, we can observe it does not give optimum results. Since the articles that are from different domains have their own way of representation of textual structure it causes a problem to train a general algorithm which will give the best results when dealt with different sorts of domains. So, in this paper we are proposing the solution regarding detection of fake news problem where we are going to use machine learning ensemble method. On Studying we have come across towards different properties which we are going to use in order to distinguish between fake news and real news. With help of these properties, we are going to train a mixture of different Machine Learning Algorithms through varied methods This has been proved to be helpful when we have to choose a large choice of applications and by implementing famous techniques or methods like bagging and boosting because these learning models can reduce or decrease the error rate to a minimum level. These techniques or methods proofs to be easier and more efficient during the training of different machine learning algorithms.

2. Literature review:

2.1 Fake News Impact:

In today's world, internet is driven by news and advertisement. Websites which contains hot news and sensational headlines results in advertising helps in capitalizing the high traffic to the site. It has been seen that fake news websites makers made money with the help of automated advertising that rewards them high traffic to their websites. The continuous spreading of information causes stress and confusion among the public or genearl citizens. fake news that deliberately created to cause damage and to mislead the general public is called or known as digital misinformation. Misinformation has a very high potential to cause problems, among minutes, for variant of people. Misinformation has been famed to cause disputes, to disrupt elections, to produce unease, and the most important hostility among the general public

2.2 Fake News influenced by social platforms:

I today's world or society, the internet became a significant or very important part of our society or in simple words of our daily lives. Anyone can easily receive the modified form of Information so easily through Social Media such as Whatsapp, Facebook etc. It was also reported in 2016 that Whatsapp and Facebook is the a very big or giant social media platform, which consists more than 20 lakh users around the world. In spreading of the fake news it has been seen that the role played by Facebook & whatsapp has the greatest impact as compared to all the other social media platforms such as twitter, orkut,etc. It had also been reported that approximately half of globe users get their news from Facebook and 22% from whatsapp. 25% of Facebook users have indicated that they have shared the false information, in many forms which also include accidental form. The spread of such news is mainly with the help of social media platforms and it's happening at a very very fast speed or pace.

Sub-Categories of the Spreading of Fake News

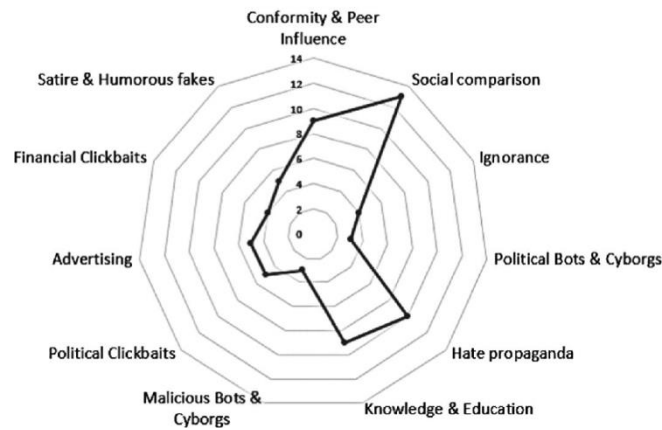


Figure 1

3. Formulation Of Problem:

In this paper (research paper), we are required to study and research about the fake news detection which includes the creators, articles, and problem subjects in different or various online social media forms and platforms. Based on various forms of this diversified data sources, therefore the article subject and various authorship & the most important relationships among them, we are required to aim for different different fake news which are generating from social networks continuously. We are required to find out a way to identify the fake news for getting good quality content, where only the important and genuine news can exist whereas fake ones either identified or be removed or scraped from social media platform. Dealing with fake news detection isn't a very easy task or a cake walk because of the following reasons:

Use of Textual data: For the articles related to news, creators of those articles and subjects, a collection or group of the information in terms of text regarding their created contents such as news, advertisement, description and profile will be collected from social media such as facebook, whatsapp, twitter and various other platforms. To perform all these actions we require a model with best feature extraction and learning ability qualities.

Fusion of Heterogeneous Information: The labels of credibility of creators and news articles have a very very strong connections, which we can indicate by or with the help of the article subject & authorship relationships among them. An efficient and very effective way of such correlations in our model learning will be very helpful for so much precise results regarding fake news.

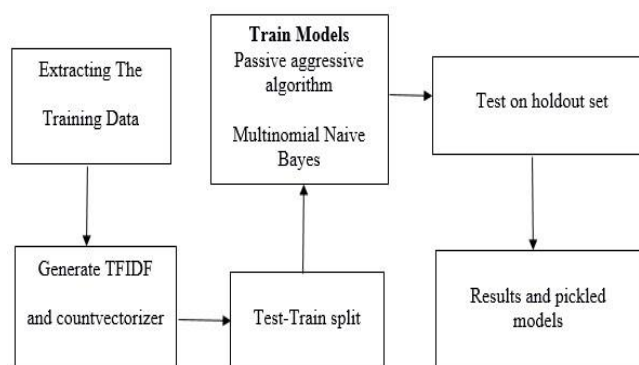


Figure 2

4. Software Requirement Analysis

- PYTHON
- NUMPY
- MATPLOTLIB
- PANDAS
- Scikit-learn

5. Feasibility Analysis

Model Performance: Our designed model is based on classification and is aimed to produce efficient results.

- Technological considerations: The analysis will be performed on a large set of data and from that only reliable sources are taken into consideration.

- Financial feasibility: The model is less expensive as we gather information from government sites which are free to access and we get the structured data from kaggle also. A large staff is also not required as the software only requires basic concepts to work on.

- Resource feasibility: The model is primarily depended on large data sets. So, having large resources will maximize the result more effective.

6. Training Algorithm, Classification of new Article workflow:

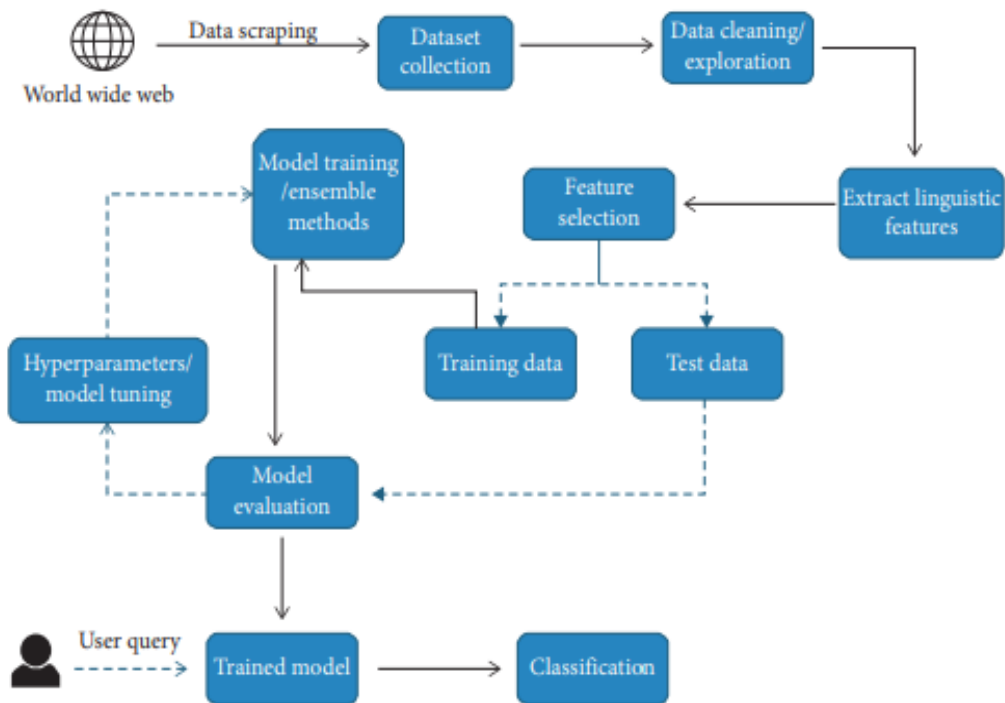


Figure 3

7. Training, Implementation and Testing

7.1 Pre processing of the given data:

In pre-processing step our first step is to import libraries then load the dataset which is open source dataset and we got it from kaggle in Jupyter Notebook. From the dataset extract the dependent and independent variables. By importing the train test split the process of splitting of dataset into training and test set occurs. After performing this process with code pre-processing of data completes.

7.2 Fitting of the Decision tree algorithm to the Training set:

In this process we fit the proposed model to the training set. For implementing this firstly import the Decision Tree Classifier class from sklearn.tree library.

7.3 Predicting the test-set result:

In this we will try to predict the test-set result by fitting the model by creating a new prediction vector.

7.4 Confusion matrix:

After getting the output we might get some incorrect predictions so for this we have to know the number of correct and incorrect predictions. For this we need the confusion matrix and for this we import the confusion matrix from sklearn.metrics.

7.5 Visualizing the training and test-set result:

In the visualization process what we will do is, we will try to visualize the test and training set result .To do this we have to plot a graph for the decision tree classifier.

8. Proposed system Merits(Advantages):

Accuracy achieved by the different methods are:

	Tf-idf vectorizer		Count vectorizer	
	Text	Title	Text	Title
Multinomial Naive Bayes	85.03	82.4	87.23	82.4
Passive Aggressive Algorithm	88.9	78.4	92.5	89.06

Figure 4

From the above accuracy table it is found that choosing count vectorizer and implementing the Passive Aggressive Algorithm on the Text data results in more accuracy.

Accuracy:

Accuracy is simply the metric which is representing, proportions of properly and precisely predicted observations or results. To calculate the accuracy of the model’s performance, the below equation is often used or mostly used:Most of the time, high accuracy is the good and more efficient and effective model.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

TP=True Positive

FP=False Positive

TN=True Negative

FN=False Negative

9. Conclusion:

The method of identifying the news manually needs a very good information of this domain and good experience to spot anomalies or errors in the given context. In this analysis, we are required to mention the matter of classifying or identifying the fake news articles with the help of machine learning model .The data we tend to employed in our work,we get it from online open source platform and it have the articles of news from a very large number of domains which covers the maximum amount of the news and also covers various domains such as political and sports news. The main idea of this research paper is identifying the pattern in the given information in text that is able to find the difference between the fake article news from the true news . detection of Fake news has several problems that needs attention of developers, data scientists, scholars and researchers. for example,in order the to reduce the fake news from spreading, important key components or steps concerned within the spread of article or story is a very important step in the way. Machine learning models and Graph theory are usually used to identify the key sources that are involved in the spread of fake articles and fake news.

References

1. IDEA:-<https://olympus.greatlearning.in/courses/14365>
2. Studied from:- <https://www.javatpoint.in/machine-learning-decision-tree-classification-algorithm>
3. <https://www.researchgate.net/publication/339022255>
4. S. B. Parikh, V. Patil, and P. K. Atrey, “On the Origin, Proliferation and Tone of Fake News,” Proc. - 2nd Int. Conf. Multimed. Inf. Process. Retrieval, MIPR 2019, pp. 135–140, 2019