IJAIR
IMPACT FACTOR
7.36
International Journal of Advance and Innovative Research
ISSN : 2394-7780

# International Journal of
# Advance and Innovative Research

## AN ANALYSIS OF DETECTING FAKE NEWS WITH PYTHON AND MACHINE LEARNING

**Dr. Sujatha Sundar Iyer[1] and Rajmohan Yadav[2]**

Assistant Professor, Department of CS/IT, Satish Pradhan Dnyanasadhana College, Thane

**ABSTRACT**

*Today organizations are spending more on newer technologies like Artificial Intelligence, Machine Learning and Deep Learning to get insight of data to perform real-world tasks and give solutions. We can call it data-driven decisions taken by machines. Nowadays, Machine Learning is the main field of computer science. It can provide sense to data. In the same way human beings can do in simple words. ML is a type of artificial intelligence. Usually, it extracts patterns out of raw data by using an algorithm. Python has a very powerful set of packages. The packages like numpy, scipy, pandas, scikit-learn etc. which are very important for machine learning and data science. Here, using python and machine learning, we are going to find a type of yellow journalism. Usually that kind of fake news is generally spread through social media. Here a model is built to find and classify news as FAKE or REAL. sklearn, TfidfVectorizer, PassiveAggressive are used in this model. The accuracy score informed us about the efficiency of the model. Here the political data set is considered. Accuracy of this model is 90% above.*

*Keywords:* Fake news, Machine learning, Linguistics, Semantics, Syntax, Algorithms, Digital tools, social media

## INTRODUCTION

World is evolving quickly. Most likely we have various benefits of this advanced world, however it has its impediments also. There are various issues in this advanced world. One of them is fake information. Somebody can without much of a stretch spread fake news. Counterfeit words are gotten out to hurt the standing of an individual or an association. It very well may be a publicity against somebody that can be an ideological group or an association. There are different web-based stages where the individual can get out the fake word. This incorporates the Facebook, Twitter and so forth AI is the piece of man-made reasoning that aids in creating the frameworks that can learn and perform various activities (Donepudi, 2019). An assortment of AI calculations is accessible that incorporate the regulated, solo, support AI calculations. The calculations initially must be prepared with an informational collection called train informational index. After the preparation, these calculations can be utilized to perform various assignments. AI is utilized in various areas to perform various undertakings.

More often than not AI calculations are utilized for expectation reasons or to distinguish something stowed away. Online stages are useful for the clients since they can without much of a stretch access some news. However, the issue is these offers the chance to the digital lawbreakers to get out a phony word through these stages. This news can be demonstrated to be hurtful to an individual or society. Per users read the news and begin trusting it without its check. Identifying the phony news is a major test since it's anything but a simple undertaking (Shu et al., 2017). In the event that the phony news isn't recognized early, then, at that point, individuals can spread it to other people and every one individual will begin trusting it. People, associations, or ideological groups can be impacted through the fake news. Individuals' suppositions and their choices are impacted by the fake news in the US appointment of 2016 (Dewey, 2016). Various analysts are working for the discovery of fake news.

The utilization of AI is demonstrating usefulness in such a manner. Scientists are utilizing various calculations to recognize the bogus news. Specialists in (Wang, 2017) said that fake news discovery is a large test. They have utilized the AI for identifying counterfeit news. Analysts of (Zhou et al., 2019) observed that the fake news is expanding with the progression of time. To that end there is a need to recognize news. The calculations of AI are prepared to satisfy this reason. AI calculations will distinguish the phony news naturally whenever they have prepared.

### The Evolution of Fake News and Fake News Detection

This is not new. Before the period of computerized innovation, it was spread through basically sensationalist reporting centered around shocking news like wrongdoing, tattle, calamities, and ironical news (Stein-Smith 2017). The commonness of fake news connects with the accessibility of broad communications computerized devices (Schade 2019). Since anybody can distribute articles by means of computerized media stages, online news stories incorporate well-informed pieces yet in addition assessment-based contentions or just bogus data (Burkhardt 2017). There is no overseer of validity norms for data on these stages making the spread of fake

news conceivable. To compound the situation, it is in no way, shape or form clear differentiating between genuine news and semi-valid or bogus news (Pérez-Rosas et al. 2018).

The idea of web-based media makes it simple to get out counterfeit words, as a client possibly sends counterfeit news stories to companions, who then, at that point, send it again to their companions, etc. Remarks on counterfeit news here and there fuel its 'believability' which can prompt fast sharing bringing about additional phony news (Albright 2017).

Social bots are likewise answerable for the spreading of phony news. Bots are here and there used to target super-clients by adding answers and notices to posts. People are controlled through these activities to share the fake news stories (Shao et al. 2018).

Misleading content is one more instrument empowering the spread of phony news. Misleading content is a publicizing instrument used to stand out enough to be noticed by clients. Shocking features or news are regularly utilized as misleading content that explore the client to notices. More taps on the advert implies more cash (Chen et al. 2015a).

Luckily, instruments have been produced for recognizing counterfeit news. For instance, an instrument has been created to recognize counterfeit words that get out through web-based media through analyzing lexical decisions that show up in features and other exceptional language structures (Chen et al. 2015b). Another apparatus, created to recognize counterfeit news on Twitter, has a part considered the Twitter Crawler which gathers and stores tweets in a data set (Atodiresei et al. 2018). Whenever a Twitter client needs to check the exactness of the news observed they can duplicate a connection into this application after which the connection will be handled for counterfeit news identification. This interaction is based on a calculation called the NER (Named Substance Acknowledgment) (Atodiresei et al. 2018).

There are numerous accessible ways to deal with assisting general society to recognize counterfeit news and this paper expects to improve comprehension of these by sorting these methodologies as found in existing writing.
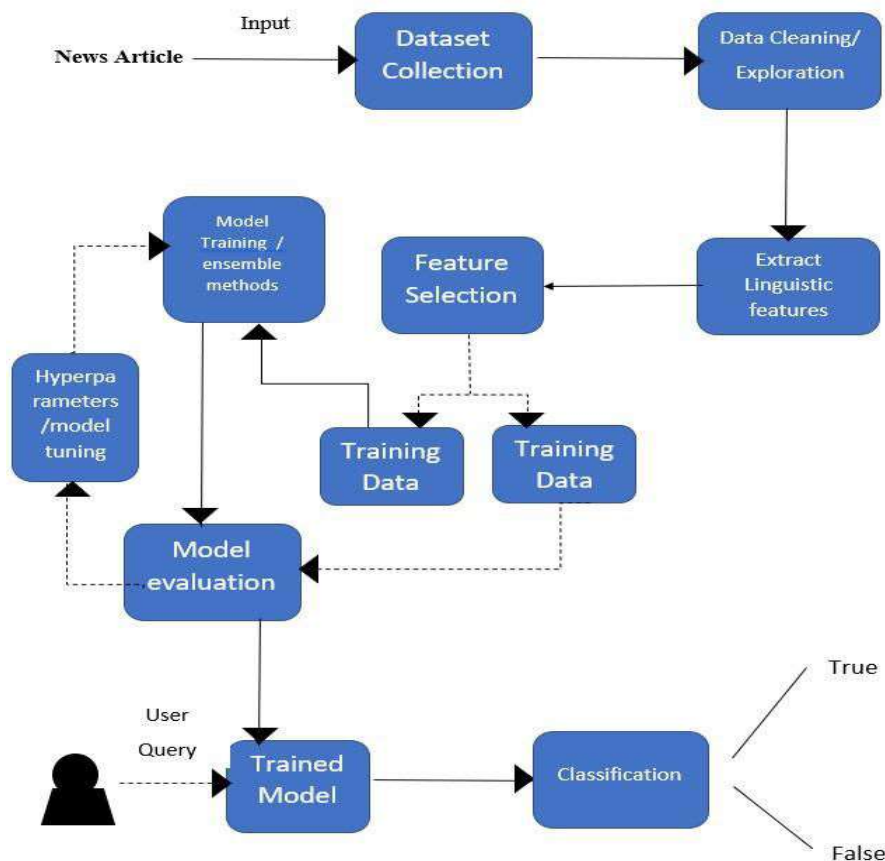
## LITERATURE REVIEW
**Shloka Gilda**, the author introduced the conception of the significance of NLP in stumbling across incorrect information. They have used TF-IDF of bigrams and probabilistic context-free grammar detection. Shloka Gilda introduced the concept of the importance of NLP in stumbling over incorrect information and examined the data set in more than one class of algorithm to find a better model which identifies noncredible resources with an accuracy of 71.2%.

**Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang and Huan Liu** , here they detect fake news on social media, which includes psychology and social theories. This article appears at principal elements answerable for the sizable attractiveness of fake messages by the user which is naive realism and confirmatory bias. Two methods used are 1) feature extraction and 2) modelling, analysing data sets, and confusion matrix for detecting fake news.

**Shivam B. Parikh and Pradeep K. Atrey**, Social networking sites read news substantially 3 ways: The (multilingual) text is analysed with the help of computational linguistics, which semantically and totally focus on text. Since utmost publications are in the form of text, a lot of work has been done on analysing them. Multimedia: Several forms of media are integrated into a single post. It may be audio, video, images, and graphics. This is very attractive and attracts the viewer's attention. Hyperlinks allow the author of the post to refer to various sources and thus gain the trust of viewers.

**Mykhailo Granik and Volodymyr Mesyura**, this paper described an artificial intelligence algorithm called the Naive Bayes classifier. The main objective of this paper is to examine how this particular method works for the particular problem.

**PROPOSED MODEL:**



The goal here is to identify whether a "news" article is fake or fact. We will take a dataset of labelled public-messages and apply classification techniques with frequency vectorizer. We can later test the model for accuracy and performance on unclassified public-messages. Similar techniques can be applied to other NLP applications like sentiment analysis etc.

We are using dataset from kaggle.com which contains the following features:

- id: unique id for a news article

- title: the title of a news article

- author: author of the news article

- text: the text of the article; could be incomplete

- label: a label that marks the article as potentially unreliable
  1: unreliable
  0: reliable

We use TfIdf Vectorizer to convert our text strings to numerical representations and initialize a PassiveAgressive Classifier to fit the model. In the end, the accuracy score and confusion matrix tell us how well our model works.

**Term Frequency(Tf) — Inverse Document Frequency(Idf) Vectorizer**
Tf-Idf Vectorizer is a common algorithm to transform text into meaningful representation of numbers. It is used to extract features from text strings based on occurrence.

We assume that a higher number of repetitions of a word would mean greater importance in the given text. We normalize the occurrence of the word with the size of the document and hence call it term-frequency. Numerical definition: tf(w) = doc.count(w) / total words in the doc

While computing term-frequency, each term is given equal weightage. There may be words which have high occurrence across the documents and hence would contribute less in deriving the meaning of the document.

Such words for example 'a', 'the' etc. might suppress the weights of more meaningful words. To reduce this effect, Tf is discounted by a factor called inverse document frequency. idf(w) = log(total_number_of_documents / number_of-documents_containing_word_w)

Tf-Idf is then computed by taking a product of Tf and Idf. More important words would get a higher tf-idf score. tf-idf(w) = tf(w) * idf(w)

## PASSIVE AGGRESSIVE CLASSIFIER

The passive-aggressive algorithms are a family of algorithms for large-scale learning. Intuitively, passive signifies that if the classification is correct, we should keep the model, and, aggressive signifies that if the classification is incorrect, update the model to adjust to more misclassified examples. Unlike most others, it does not converge, rather it makes updates to correct the loss.

## DEVELOPING THE ML MODEL

**Step 1:** Import the necessary packages:

```
import numpy as np
import pandas as pd
import itertools
import seaborn as sn
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import PassiveAggressiveClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

**Step 2:** Load the dataset into pandas' data-frame:

```
train = pd.read_csv('train.csv')
test = pd.read_csv('test.csv')
test = test.set_index('id', drop = True)
```

**Step 3:** Read and understand the data. One of the most import steps while creating any ML model is to first prepare the data. This includes cleaning and filtering the data, removing outliers, and creating feature that are independent and sensible

We use. shape method to identify number of columns in the dataset and the total number of news samples. Then read the data table using. head method to see how the data looks. Next, identify column names where news articles are written and the ones where classification is marked.

We then use. isna to identify if we have any null values in the column where our news articles are put, in this case it is in the column named 'text'. Now, we use. sum () to identify how many such values exist. Once identified, we drop the rows where the column 'text' has null values and fill a blank space in other columns with null values.

```
# Counting number of rows and columns in the data
print('Shape of Training Data: ', train.shape)
```

```
# Gettiing a hang of the data in each column and their names
print('\n \n TRAIN \n', train.head())
print('\n \n TEST \n', test.head())
```

```
# Looking for any places where training data has NaN values
print('\n \nNumber of Null values in Train Set: ', train['text'].isna().sum())
print('Number of Null values in Test Set: ', test['text'].isna().sum())
```

```
# Dropping all rows where text column is NaN
train.dropna(axis=0, how="any", thresh=None, subset=['text'], inplace=True)
test = test.fillna(' ')
```

**Step 4:** Let us now see if we have any outliers in the data. We will do this by checking the number of words in each article and identifying the range and mean of the number of words in all articles. We will use len() function to check for the lengths.

**Step 5:** One final step before we start applying the model is to segregate the classification column with the rest of the input features, and then dividing the dataset into training and testing subsets. We do this split to ensure that our model performs well on a new dataset. We take 90% of our data as the training set and 10% as the testing set. This split percentage can be customized in order to tune the model better.

**Step 6:** Let's initialize a TfIdfVectorizer with stop words from the English language and a maximum document frequency of 0.7 (terms with a higher document frequency will be discarded). Stop words are the most common words in a language that are to be filtered out before processing the natural language data. And a TfIdfVectorizer turns a collection of raw documents into a matrix of Tf-Idf features.

**Step 7:** Next, we'll initialize a PassiveAggressiveClassifier. We'll fit this on tfidf_train and y_train.

## RESULT
In fake news detection, supervised and unsupervised learning algorithms are used to classify text. In this review paper, we try to find the solution for the fake news detection problem using the machine learning approach. We observed that the Random Forests algorithm with a simple term frequency-inverse document frequency vector gives the best output compared to others. Our study examines various text properties that can be used to distinguish fake and real content.

## ADVANTAGES
Fake News Detection system will help in controlling the spread of fake news over social media. This way, we can help the people to make more informed decisions, and they are not made to think about what others are trying to manipulate to believe. A Fake News Detection system will reduce the burden to check the authenticity of the news manually and saves lots of time.

## DISADVANTAGES
The accuracy of detecting fake news will not be 100%. Therefore, some articles may be predicted as false.

## CONCLUSION
Nowadays, more people are constantly consuming news from social media in place of the conventional media. This fake news develops a sturdy bad effect on users and the society. Therefore, for detecting the fake news, examine specific studies and identify Word Embedding, Tokenization and Parts of speech tagging are best for Pre-Processing of data and also identifies TF-IDF and Count Vectorizer are best for feature extraction. So, further we want to use those methods for PreProcessing, feature extraction and also, we want to implement the Random Forest classifier, Convolutional Neural Networks, Long Short-Term Memory for high accuracy and an Ensemble Learning Approach for high accuracy. It takes a lot of time to verify a single article manually. That's why we have discussed the problem of classifying fake news articles using machine learning models. This way, we can help the people to make more informed decisions, and they won't be led to think about what others are trying to manipulate them into believing

## REFERENCE
Abdullah-All-Tanvir, Mahir, E. M., Akhter S., & Huq, M. R. (2019). Detecting Fake News using Machine Learning and Deep Learning Algorithms. 7th International Conference on Smart Computing & Communications (ICSCC), Sarawak, Malaysia, Malaysia, 2019, pp.1-5, https://doi.org/10.1109/ICSCC.2019.8843 612

Ahmed, H., Traore, I., & Saad, S. (2017). Detection of online fake news using n-gram analysis and machine learning techniques. Proceedings of the International Conference on Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments, 127–138, Springer, Vancouver, Canada, 2017. https://doi.org/10.1007/978-3-319-69155- 8_9

Ahmed, H., Traoré, I., & Saad, S. (2018). Detecting opinion spams and fake news using text classification. Secur. Priv., 1(1), 1-15. https://doi.org/10.1002/spy2.9

Al Asaad, B., & Erascu, M. (2018). A Tool for Fake News Detection. 2018 20th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC), Timisoara, Romania, 2018, pp.379-386. https://doi.org/10.1109/SYNASC.2018.00 064

Aphiwongsophon, S., & Chongstitvatana, P. (2018). Detecting Fake News with Machine Learning Method. 2018 15th International Conference on Electrical Engineering/Electronics, Computer, Telecommunications and Information Technology (ECTI-CON), 528-531. https://doi.org/10.1109/ECTICon.2018.86 20051

Della Vedova, M. L., Tacchini, E., Moret, S., Ballarin, G., DiPierro, M., & de Alfaro, L. (2018). Automatic online fake news detection combining content and social signals. FRUCT'22: Proceedings of the 22nd Conference of Open Innovations Association FRUCT. Pages 272–279. https://dl.acm.org/doi/10.5555/3266365.32 66403

Dewey, C. (2016). Facebook has repeatedly trended fake news since firing its human editors. The Washington Post, Oct. 12, 2016.

Donepudi, P. K. (2019). Automation and Machine Learning in Transforming the Financial Industry. Asian Business Review, 9(3), 129-138. https://doi.org/10.18034/abr.v9i3.494

Donepudi, P. K. (2020). Crowdsourced Software Testing: A Timely Opportunity. Engineering International, 8(1), 25- 30. https://doi.org/10.18034/ei.v8i1.491

Donepudi, P. K., Ahmed, A. A. A., Saha, S. (2020a). Emerging Market Economy (EME) and Artificial Intelligence (AI): Consequences for the Future of Jobs. Palarch's Journal of Archaeology of Egypt/Egyptology, 17(6), 5562- 5574. https://archives.palarch.nl/index.php /jae/article/view/1829 11) Donepudi, P. K., Banu, M. H., Khan, W., Neogy, T. K., Asadullah, ABM., Ahmed, A. A. A. (2020b). Artificial Intelligence and Machine Learning in Treasury Management: A Systematic Literature Review. International Journal of Management, 11(11), 13-22.

Jadhav, S. S., & Thepade, S. D. (2019). Fake News Identification and Classification Using DSSM and Improved Recurrent Neural Network Classifier, Applied Artificial Intelligence, 33(12), 1058- 1068, https://doi.org/10.1080/08839514.20 19.1661579

Kaliyar, R. K., Goswami, A., Narang, P., & Sinha, S. (2020). FNDNet–A deep convolutional neural network for fake news detection. Cognitive Systems Research, 61, 32-44. https://doi.org/10.1016/j.cogsys.2019.12.0 05

Kaur, S., Kumar, P. & Kumaraguru, P. (2020). Automating fake news detection system using multi-level voting model. Soft Computing, 24(12), 9049–9069. https://doi.org/10.1007/s00500-019- 04436-y

Kesarwani, A., Chauhan, S. S., & Nair, A. R. (2020). Fake News Detection on social media using K-Nearest Neighbor Classifier. 2020 International Conference on Advances in Computing and Communication Engineering (ICACCE), Las Vegas, NV, USA, pp.1-4, https://doi.org/10.1109/ICACCE49060.20 20.9154997