# Using Different machines learning model to detect fake news.

*Abstract*—**Fake news is seditious propaganda aimed at spreading false information among the population through news broadcasts or social networks. This design will help determine the delicacy of fake news and real news. We've two different datasets. To preprocess our dataset, we went through five important way. After executing the algorithm, we gain the results of each of the 3 classifiers. From our algorithm results, we can see that support vector class has the topmost delicacy of 100.**

*Index Terms*—**fake news, detection, accuracy, support vector classification**

## I. INTRODUCTION

Fake news is seditious propaganda aimed at spreading false information among the population through news broadcasts or social networks. It's worth noting that fake news spreads quick than real news. Researchers at the Massachusetts Institute of Technology claim that false news disseminates ten times more quickly than legitimate information and is defined as claims that have been debunked by six prominent fact-checking websites. While several US senators and other opponents decried the dissemination of incorrect information on automated bots before the 2016 election, MIT experimenters vetted tweets generated by bots used for their inquiry. Still, the perpetration of artificial intelligence has effectively paved the way for stopping the spread of fake news. AI trains machines grounded on their capability to achieve asked pretensions. It introduces different models to descry intimidating news with fake and real status. These models actually help reduce the lifetime of fake news. This design will help descry the delicacy of fake news and real news through collecting different types of data by enforcing certain machine literacy models. This will take stoner feedback on all news, dissect it and better the affair by informing about real or fake news.

## II. METHODOLOGY

### A. Dataset Description

The datasets we used were preliminarily used for two reviews scratched by H Ahmed,I Traoré and S Saad [2] [3]. Two distinct types of datasets are basically used to descry fakes and genuine goods. One set of data is used as fake and the other as real. 17,903 unique values are represented in the combined fake or real news dataset. further than 2000 pieces of data are contained in the dataset to fantasize true or false using a four- order model.2.2.

### B. Dataset Restructure

We've two different datasets and before incorporating them, we produce a new column for each dataset called " Marker ". Now we label all Fake News = 0 and Real News = 1. also we combine the two datasets into one dataset, mix them and do the rest of the work with this data set.

### C. Pre-Processing Techniques Applied

To preprocess our dataset, we went through five important way.

*Null Checking:*
The value null indicates that a variable doesn't point to any object and contains no value. It's frequently used to specify or authenticate the virtuality of commodity. Then we've checked to see if there's any missing data using a introductory " if " statement. Since we did not find any null values, deleting the rows isn't necessary.

*Produce" Content" content :*
It principally involves incorporating all the information from an composition into a single textbook. It contains Title, Subject and Text values. Remove gratuitous columns This means reducing the number of columns without affecting crucial information.

*Dropping Unnecessary Columns:*
Title, subject, and textbook are integrated into the content column. This way, they no longer need to enthrall separate columns in the data set. Also, the date column values are not authentically important for our design. So we removed these columns. And keep only the necessary columns Content and Marker.

*Stemming:*
Stemming is the process of reducing a word's declension to its root form. As an example, a set of words may be mapped to the same stem even if the stem itself may not be a recognized term in the language.This sort of word normalization is used.This is a method of docking the hunt by turning a string of words from a judgment into a set. words that, although having the same meaning, have a different meaning in some contexts or when used in formal contexts. By regaining access to stoner-entered expressions, additional documents are matched since stoner-entered expressions' necessary word forms are likewise matched, increasing overall recall. This comes at the cost of reduced delicacy. still, before

stemming, we removed allnon-letter characters, converted all letters to lowercase, and also stemmed all words. And later each that, we removed the stop words.

*Vectorization:*

Vectorization is a fashion for performing array operations without using for circles [5]. Instead, we employ colorful, highly efficient modules to construct functions that shorten the length and frequency of legal proceedings. By transforming the text in the textbook into digital vectors, vectorization is utilized to highlight certain distinctive elements from the text that the model can train on.

## D. Models applied

*Logistic Regression Score:*

The logistic regression function $p(x)$ is the sigmoid function of $f(x)$ :

Consequently, it often ranges from 0 to 1. Because the logistic regression (LR) model offers the fundamental equation required for categorizing issues into two or more categories, it is employed. In reality, we divide textbooks into two categories based on a wide variety of features (true/wrong composition or correct/wrong composition). We performed hyperparameter optimization to acquire the official results for each distinct dataset.

$$h_\theta(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} \quad (1)$$

The results of logistic regression are converted into probability values using the sigmoid function; The thing is to minimize the cost function to achieve the formal probability. When calculating the cost function shows:

Classification of support vectors : The support vector machine (SVM), which contains a wide range of kernel functions, is another model for the double classification issue. The goal of the SVM model is to categorize data points by computing a hyperplane (or decision boundary) based on a collection of characteristics.(7) The objective is to identify the hyperplane dividing the data points of the two classes with the biggest perimeter in a $mathrm$-dimensional space, which can take on a wide variety of geometries in an N-dimensional environment. The SVM model's cost function is well illustrated:

$$J(\theta) = \frac{1}{2} \sum_{j=1}^{n} \theta_j^2 \quad (2)$$

$$\theta^T x^{(i)} \geq 1, \quad y^{(i)} = 1 \quad (3)$$

$$\theta^T x^{(i)} \leq -1, \quad y^{(i)} = 0 \quad (4)$$

The law that follows employs a direct kernel. Kernels are often used to fit data points that are multidimensional or challenging to separate cleanly. The sigmoid SVM, kernel SVM (polynomial SVM), Gaussian SVM, and fundamental direct SVM models were all employed in our script.

*Naive Bayes classifier:*

This classification strategy uses the Bayes theorem and the predictor independence presumption (b8). The Naive Bayes classifier, to put it simply, assumes that every point in a class exists independently of every other point. This approach for categorizing objects is appropriate for both double and multiclass classification. Because the computation of chances for each class is streamlined to make their computation simpler, it is also known as Diot Bayes. When input variables are categorical as opposed to numeric, Naive Bayes performs well. Making data predictions and forecasts based on concrete facts is advantageous.

## 3 Results

This is Logistic Regression's Result:

| Logistic Regression | | | | |
|---|---|---|---|---|
| | precision | Recall | F1 score | Support |
| 0 | 0.45 | 0.45 | 0.45 | 4384 |
| 1 | 0.45 | 0.45 | 0.45 | 4041 |
| avg | 0.45 | 0.45 | 0.45 | 8425 |

This is SVM's Result:

| SVM | | | | |
|---|---|---|---|---|
| | precision | Recall | F1 score | Support |
| 0 | 0.42 | 0.42 | 0.42 | 4384 |
| 1 | 0.42 | 0.42 | 0.42 | 4041 |
| avg | 0.42 | 0.42 | 0.42 | 8425 |

This is Naive Bayes's Result:

| Naive Bayes Result | | | | |
|---|---|---|---|---|
| | precision | Recall | F1 score | Support |
| 0 | 1.00 | 1.00 | 1.00 | 4682 |
| 1 | 1.00 | 1.00 | 1.00 | 4298 |
| avg | 1.00 | 1.00 | 1.00 | 8425 |

After enforcing one of the test data using our Machine Learning algorithm. The algorithm presents us with the results it discovered by learning from the training data we preliminarily stored. Data mining has instructed us that papers like " No author, No title, No image and bs " are more likely to be fake news. In all classification cases, we resolve the dataset into 20 testing data and 80 training data.

| Classifiers | precision | Recall | F1 Score |
|---|---|---|---|
| LR | 0.45 | 0.45 | 0.45 |
| SVM | 0.42 | 0.42 | 0.42 |
| NB | 1.00 | 1.00 | 1.00 |

The model's positive vaticination values( perfection) represent applicable textbook among the recaptured textbook documents, while perceptivity( recall) represents the bit of the total textbook documents related has actually been recaptured. The F1 score represents the harmonious mean of the combination of perfection and recall. Classifier Precision recall f1- Score

After enforcing the algorithm, we attained the perfection, recall, and f1 score of each of the three classifiers. From our algorithm results, we can see that support vector classification has the topmost delicacy of 100%.

## REFERENCES

[1] "Fake News Spreads Fast, but Don't Blame the Bots." Internet Society, 21 Mar. 2018, https:// www.internetsociety.org/blog/2018/03/

[2] Ahmed H, Traore I, Saad S. "Detecting opinion spams and fake news using text classification", Journal of Security and Privacy, Volume 1, Issue 1, Wiley, January/February 2018.

[3] Ahmed H, Traore I, Saad S. (2017) "Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques. In: Traore I., Woungang I., Awad A. (eds) Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments. ISDDC 2017. Lecture Notes in Computer Science, vol 10618. Springer, Cham (pp. 127-138).

[4] "Stemming and Lemmatization in Python." DataCamp, DataCamp, 23 Oct. 2018, https://www.datacamp.com/tutorial/ Stemming-lemmatization-python.

[5] "Vectorization in Python - A Complete Guide." AskPython, 5 Sept. 2021, https: //www.askpython.com/python-modules/ numpy/vectorization-numpy.

[6] Real Python. "Logistic Regression in Python." Real Python, Real Python, 18 Aug. 2022, https://realpython.com/ logistic-regression-python/

[7] Contributor, TechTarget. "What Is Support Vector Machine (SVM)? - Definition from Whatis.com." WhatIs.com, TechTarget, 29 Nov. 2017, https://www.techtarget.com/whatis/

[8] Brownlee, Jason. "Naive Bayes Classifier from Scratch in Python." Machine Learning Mastery, 24 Oct. 2019, https://machinelearningmastery.com/