

ABSTRACT :- Social media is one of the very powerful media in spreading information. People are interested in sharing without any proper checking of any sort of false information. Unstructured text data may be classified into meaningful categorical classifications using text classification, which is a typical study area in the discipline of Natural Language Processing (NLP). The main contribution of this article is to identify a finest framework to tackle the fake news problem with the NLP and Machine Learning techniques. In this empirical research, the fake news data is analysed with the different combinations of Vectorizers and Machine Learning Classifiers. From the experimental results on five benchmark datasets namely fake_real_news dataset extracted from Kaggle, COVID-19 Constrain, Politifact, ISOT and Gossipcop, it is observed that the fake news detection with the combination of TF-IDF Vectorizer and Passive-Aggressive Classifier outperforms the other existing methods. **Keywords:** Machine Learning, Pre-processing, Vectorizer, Classifiers, Natural Language Processing.

The term “fake news” refers to a tale that intentionally misleads the public, but in recent months, social media has begun to redefine the term. The importance of disinformation during the time of COVID19 pandemic is a subject of weighty attention, particularly following the emergence of multiple variants of Coronavirus. Factually false and deceptive stories produced primarily for the goal of gaining money through page views became known as “fake news.” With the use of this research, we want to develop a model that can properly identify whether or not a piece of content is fake news. WhatsApp in India has become one of the Major fakenews providers.

METHODOLOGY

Feature Extraction Vectorizers like TF-IDF and Count Vectorizer are used to do the word embeddings and dataset cleaning in this approach. It's done by removing stop words like “the,” “when,” and “there,” as well as utilising an n-number of the most often used words, phrases, lowercased or not, as well as using keywords that appear at least a certain number of times in a certain text corpus. The dataset is then divided into two parts: data for testing and data for training. Classifiers like Passive Aggressive Classifier, Multinomial Nave Bayes, Random Forest, and Support Vector Machine are fed the train data after the dataset has been divided. The classifiers train the data and check the relativity of the features in the test data and gives the result as ‘Fake’ or ‘Real’ along with the accuracy of the Classifiers.

Fake News

A sort of sensationalist reporting, counterfeit news embodies bits of information that might be lies and is, for the most part, spread through web-based media and other online media. This is regularly done to further or force certain kinds of thoughts or for false promotion of products and is frequently accomplished with political plans . Such news things may containbogus and additionally misrepresented cases and may wind up being virtualized by calculations, and clients may wind up in a channel bubble.

Data Analysis

Here I will explain the dataset.

In this python project, we have used the CSV dataset. The dataset contains 7796 rows and 4 columns .

This dataset has four columns,

- 1 . title: this represents the title of the news .
- 2 . author: this represents the name of the author who has written the news .
- 3 . text: this column has the news itself.
- 4 . label : this is a binary column representing if the news is fake (1) or real (0).

Libraries

The very basic data science libraries are sklearn, pandas, NumPy e.t.c and some specific libraries such as transformers .

```
import pandas as pd
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
import re
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import HashingVectorizer
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer

# Read dataset from CSV File
df = pd.read_csv('fake-news/train.csv')
df.head()
```

output:-

| df.head() | | | | |
|--|----|---|--------------------|---|
| executed in 16ms, finished 09:37:16 2021-06-07 | | | | |
| | id | title | author | text label |
| 0 | 0 | House Dem Aide: We Didn't Even See Comey's Let... | Darrell Lucas | House Dem Aide: We Didn't Even See Comey's Let... 1 |
| 1 | 1 | FLYNN: Hillary Clinton, Big Woman on Campus - ... | Daniel J. Flynn | Ever get the feeling your life circles the rou... 0 |
| 2 | 2 | Why the Truth Might Get You Fired | Consortiumnews.com | Why the Truth Might Get You Fired October 29, ... 1 |
| 3 | 3 | 15 Civilians Killed In Single US Airstrike Hav... | Jessica Purkiss | Videos 15 Civilians Killed In Single US Aistr... 1 |
| 4 | 4 | Iranian woman jailed for fictional unpublished... | Howard Portnoy | Print 'nAn Iranian woman has been sentenced to... 1 |

Before proceeding, we need to check whether a null value is present in our dataset or not.

`df = df.isnull()` There is no null value in this dataset. But if you have null values present in your dataset then you can fill it. In the code given below, I will tell you how you can replace the null values .
`df = df.fillna(' ')`

Data Preprocessing

In data processing, we will focus on the text column on this data which actually contains the news part.

We will modify this text column to extract more information to make the model more predictable. To extract information from the text column, we will use a library, which we know by the name of 'nltk' .

Here we will use functionalities of the 'nltk' library named Removing Stopwords, Tokenization, and Lemmatization . So we will see these functionalities one by one with these three examples . Hope you will have a better understanding of extracting information from the text column after this .

Removing Stopwords:-

These are the words that are used in any language used to connect words or used to declare the tense of sentences . This means that if we use these words in any sentence they do not add much meaning to the

context of the sentence so even after removing the stopwords we can

Tokenization is the process of breaking text into smaller pieces which we know as tokens .

Each word, special character, or number in a sentence can be depicted as a token in NLP.

Tokenization is the process of breaking down a piece of code into smaller units called tokens .

```
from nltk.tokenize import word_tokenize
```

```
text = "Hello everyone. Welcome to Analytics Vidhya. You are studying NLP article" word_tokenize(text)
```

The output looked like this:

```
['Hello everyone.', 'Welcome to Analytics Vidhya.', 'You are studying NLP article']
```

CONVERTING LABELS:-

The dataset has a Label column whose datatype is Text Category. The Label column in the dataset is

classified into two parts, which are denoted as Fake and Real. To train the model, we need to convert the label column to a numerical one.

```
df.label = df.label.astype(str) df.label = df.label.str.strip() dict = { 'REAL' : '1' ,  
'FAKE' : '0' }
```

```
df['label'] = df['label'].map(dict)df.head()
```

To proceed further, we separate our dataset into features(x_df) and targets(y_df).

```
x_df = df['total'] y_df = df['label']
```

VECTORIZATION

Vectorization is a methodology in NLP to map words or phrases from vocabulary to a corresponding vector of real numbers which is used to find word predictions, word similarities/semantics .

To make documents' corpora more relatable for computers, they must first be converted into some numerical structure. There are few techniques that are used to achieve this such as 'Bag of Words'.

Here, we are using vectorizer objects provided by Scikit- Learn which are quite reliable right out of the box.

```
from sklearn.feature_extraction.text import TfidfTransformer
```

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
count_vectorizer = CountVectorizer()
```

```
count_vectorizer.fit_transform(x_df)
```

```
freq_term_matrix = count_vectorizer.transform(x_df)
```

```
tfidf
```

```
=
```

```
TfidfTransformer(norm
```

```
=
```

```
"l2")
```

```
tfidf.fit(freq_term_matrix)
```

```
tf_idf_matrix
```

```
=
```

```
tfidf.fit_transform(freq_term_matrix)
```

```
print(tf_idf_matrix)
```

Here, with 'Tfidftransformer' we are computing word counts using

'CountVectorizer' and then computing

the IDF values and after that the Tf- IDF scores . With 'Tfidfvectorizer' we can do all three steps at once.

The code written above will provide with you a matrix representing your text. It will be a sparse matrix with a

large number of elements in a Compressed Sparse Row format.

The most used vectorizers are:

Count Vectorizer: The most straightforward one, it counts the number of times a token shows up in the document and uses this value as its weight.

Hash Vectorizer: This one is designed to be as memory efficient as possible. Instead of storing the tokens as strings, the vectorizer applies the hashing trick to encode them as numerical indexes. The downside of this method is that once vectorized, the features' names can no longer be retrieved.

TF-IDF Vectorizer: TF- IDF stands for "term frequency-inverse document frequency", meaning the weight assigned to each token not only depends on its frequency in a document but also how recurrent that term is in the entire corpora.

Fake News Detection using NLP techniques

Fake news detection is a hot topic in the field of natural language processing. In this article, we are using this dataset for news classification using NLP techniques. We are given two input files.

One with real news and the other one with fake news.

Our job is to create a model which predicts whether a given news is real or fake. As this is a supervised learning problem, we are creating a target column named 'label' in both real and fake news data and concatenating them.

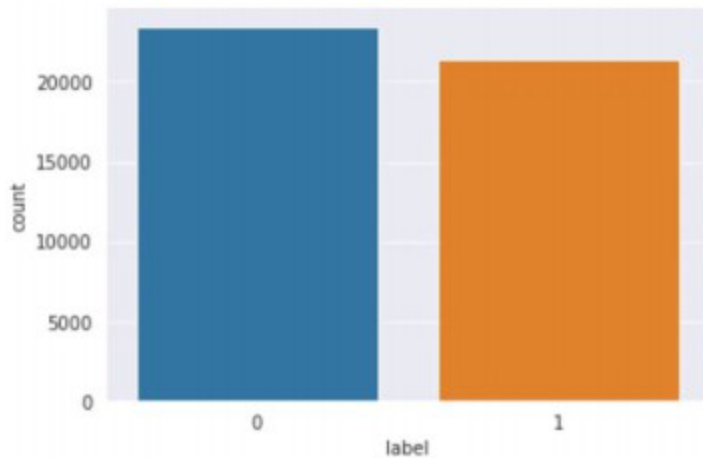
```
real['label'] = 1
```

```
fake['label'] = 0
```

```
data = pd.concat([real, fake])
```

Now we have the input where real news has the value of label as 1 and fake news have the value of label as 0. We have to check whether our data is balanced. We use seaborn library to plot the counts of real and fake news.

```
import seaborn as sns
sns.set_style("darkgrid")
sns.countplot(data['label'])
```



We can conclude from the plot that the data is balanced. This is an important step because there are a lot of real world datasets that are imbalanced. Let us check for null values in the data.

```
title      0
text       0
subject    0
date       0
label      0
dtype: int64
```

Classification

Now, let us move into the classification models. We will try different models and evaluate the performance. As classification is a supervised learning, we have to first split the data into training and test data. We train the model using the train data, and test the performance of our model using test data. Generally, the data is split in such a way that we have 80% of our data in train set, and 20% of our data in test set. This is because of the fact that the more the training set, the more the model learns from the data.

#splitting data for training and testing

```
import sklearn
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test =
train_test_split(data['text'],data['label'],test_size=0.2,
random_state = 1
```

Text Preprocessing:

Tokenization: Split the text into words or subword tokens.

Lowercasing: Convert all text to lowercase to ensure uniformity.

Removing Stop Words: Eliminate common words (e.g., "the," "and") that may not carry significant meaning.

Punctuation Removal: Strip punctuation marks from the text.

Lemmatization or Stemming: Reduce words to their base or root form.

Feature Extraction: TF-IDF (Term Frequency-Inverse Document Frequency): Calculate the importance of words in a document relative to the entire corpus.

Word Embeddings: Utilize pre-trained word embeddings like Word2Vec, GloVe, or FastText to represent words as dense vectors.

N-grams: Capture word sequences and relationships by using bigrams or trigrams.

Model Selection:

Supervised Learning: Choose a classification algorithm such as Logistic Regression, Naive Bayes, Random Forest, or deep learning models like LSTM or BERT for binary (real/fake)

classification.

Unsupervised Learning: Consider clustering techniques to identify patterns in the data that may indicate fake news.

Labeling Data: Ensure you have labeled data, typically with labels like "real" and "fake." You may also consider different classes, such as "partially true" or "satire."

Training and Testing: Split your data into training and testing sets to train and evaluate your model's performance.

Evaluation Metrics: Use metrics like accuracy, precision, recall, F1-score, and ROC-AUC to evaluate the model's performance.

Ensemble Methods: Combine multiple models or use ensemble techniques (e.g., stacking) to improve overall performance.

Explainability:

Employ techniques like LIME (Local Interpretable Model-agnostic Explanations) or SHAP (SHapley Additive exPlanations) to make model predictions more interpretable and transparent.

Feature Importance:

Analyze feature importance to understand which words or phrases contribute most to the classification decisions.

Cross-validation:

Use techniques like k-fold cross-validation to assess your model's robustness and reduce overfitting.

Real-time Monitoring:

Continuously monitor your model's performance and update it to adapt to emerging trends in fake news propagation.

User Interface:

Create a user-friendly interface or application for users to submit and check the credibility of news articles.

External Data Sources:

Incorporate external fact-checking services and databases to cross-verify information.

Metadata Analysis:

Consider analyzing metadata such as the source of the news, publication date, and social sharing data to enhance your analysis.

Language and Sentiment Analysis:

Analyze sentiment and linguistic patterns that may indicate misinformation or bias.

Community Reporting:

Enable users to report suspicious content for further investigation.

Deep Learning Architectures:

Experiment with advanced deep learning architectures like CNNs (Convolutional Neural Networks) and Transformers (e.g., BERT, GPT) to capture intricate linguistic patterns in text. Semantic Analysis: Implement semantic analysis to understand the meaning and context of text, which can help identify misleading or deceptive content.

Cross-lingual Analysis:

Extend your model to handle multiple languages for a broader scope of fake news detection.

Fact-Checking Integration:

Integrate external fact-checking organizations' APIs to cross-reference claims in news articles with their fact-checking results.

Social Network Analysis:

Analyze how news articles are shared on social media platforms, including user comments, to gauge their credibility.

Network Analysis:

Explore network-based approaches to identify sources and patterns.

Continuous Training:

Implement a continuous learning system that adapts to new types of misinformation and misinformation tactics.

Collaboration:

Collaborate with researchers, organizations, and platforms focused on combating fake news to share knowledge and resources.