

ABSTRACT

This project seeks to create a Natural Language Processing (NLP) model for distinguishing between genuine and fake news. The data used comprises two CSV files containing the corresponding news. Data pre-processing steps like tokenization, stemming, and stop word removal are incorporated. Then the data is divided into training and testing sets. TFIDF vectorization is utilized to convert text data into a numerical format. Two machine learning models, Logistic Regression and Passive Aggressive Classifier, are employed in building this classifier model. The models' accuracies are assessed by use of the scikit-learn library's accuracy score metric. This model has potential utility in combating false information spread on the internet.

OBJECTIVE

The ideal of this design is to develop a natural language processing model utilizing Python that can directly descry whether a news composition is imitative or genuine. The model will be trained on a dataset conforming of two sets of news papers one containing imitative news papers and the other containing genuine news papers. The model will pre-process the data by removing stop words, stemming the words, and vectorising the textbook utilizing the TF- IDF system. Eventually, the model will exercise two engine literacy algorithms(Logistic Regression and Passive Aggressive Classifier) to codify the news papers as imitative or genuine, and estimate its delicacy utilizing the test dataset.

INTRODUCTION

In moment's world, the spread of fake news is a major concern. With the wide use of social media platforms, anyone can partake information without vindicating its authenticity. Fake news can beget significant detriment to individualities, associations, and society as a whole. thus, the discovery of fake news is of utmost significance. In this design, we aim to use Natural Language Processing(NLP) ways to distinguish between fake and genuine news papers. We use the Python programming language and the NLTK library to pre-process the data, including tokenization, stemming, and stop word junking. We use a dataset conforming of two CSV lines containing real and fake news papers. We combine these datasets, pre-process them, and resolve them into training and testing sets. We use the TF- IDF vectorization fashion to convert the textbook data into a numerical format suitable for machine literacy algorithms. We also make two machine literacy models, videlicet Logistic Retrogression and Passive- Aggressive Classifier, to classify news papers as fake or genuine. Eventually, we estimate the delicacy of the models and choose the bone

with better performance. The proposed system can be used to automatically identify fake news papers, which can help help their spread and reduce the impact of false information on individualities and society

METHODOLOGY

The methodology for the project of natural language processing to distinguish fake news and genuine news involves the following steps: 1. Data Collection: The first step is to collect the data in the form of fake news and genuine news articles. For this project, we have used two datasets: Fake.csv and True.csv.

2. Data Pre-processing: In this step, we have performed various pre-processing steps to clean and prepare the data for analysis. This includes tokenization, stemming, stop words removal, and concatenation of the fake and genuine news articles.

3. Data Splitting: The next step is to split the dataset into training and testing sets. We have used a 75:25 split ratio for this project.

4. Vectorization: We have used the TF-IDF vectorization technique to convert the text data into numerical form, which can be used for analysis by machine learning algorithms.

5. Building of ML Models: We have used two machine learning algorithms - Logistic Regression and Passive Aggressive Classifier - to build models for distinguishing fake news from genuine news.

6. Model Evaluation: Finally, we have evaluated the performance of both models using accuracy score as the evaluation metric.

The output of the code includes the accuracy scores for both models. The Logistic Regression model achieved an accuracy of 98.9%, while the Passive Aggressive Classifier achieved an accuracy of 99.5%. These results indicate that both models are effective in distinguishing fake news from genuine news.

CODE

```
pip install nltk

import nltk

#nltk.download()

import pandas as pd

fake = pd.read_csv("Fake.csv")
genuine= pd.read_csv("True.csv")

display (fake.info())

display (genuine.info())

display(genuine.head(10))

display(fake.subject.value_counts())

fake['target']=0
genuine['target']=1

display(genuine.head(10))

display(fake.head(10))

data=pd.concat([fake,genuine],axis=0)

date=data.reset_index(drop=True)

data=data.drop(['subject','date','title'],axis=1)

print(data.columns)

"""### Tokenization"""

import nltk
nltk.download('punkt')

from nltk.tokenize import word_tokenize

data['text']=data['text'].apply(word_tokenize)

"""## Stemming"""
```

```

print(data.head(10))

from nltk.stem.snowball import SnowballStemmer
porter = SnowballStemmer("english")

def stem_it(text):
    return [porter.stem(word) for word in text]

data['text']=data['text'].apply(stem_it)

print(data.head(10))

"""### Stop word removal"""

#from nltk.corpus import stopwords

def stop_it(t):
    dt=[word for word in t if len(word)>2]
    return dt

data['text']=data['text'].apply(stop_it)

print(data.head(10))

data['text']=data['text'].apply(' '.join)

"""### Splitting"""

from sklearn.model_selection import train_test_split
X_train, X_test, y_train,
y_test=train_test_split(data['text'],data['target'],test_size=0.25)
display(X_train.head())
print('\n')
display(y_train.head())

"""### Vectoization"""

from sklearn.feature_extraction.text import TfidfVectorizer
my_tfidf = TfidfVectorizer( max_df=0.7)

tfidf_train = my_tfidf.fit_transform(X_train)
tfidf_test = my_tfidf.transform(X_test)

print(tfidf_train)

"""## Logistic Regression"""

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier

```

```
from sklearn.metrics import accuracy_score
```

```
model_1= LogisticRegression(max_iter=900)
```

```
model_1.fit(tfidf_train, y_train)
```

```
pred_1 = model_1.predict(tfidf_test)
```

```
cr1  = accuracy_score(y_test,pred_1)
```

```
print(cr1*100)
```

```
"""## PassiveAggressiveClassifier"""
```

```
from sklearn.linear_model import PassiveAggressiveClassifier
```

```
model = PassiveAggressiveClassifier(max_iter=50)
```

```
model.fit(tfidf_train, y_train)
```

```
y_pred = model.predict(tfidf_test)
```

```
accscore = accuracy_score(y_test, y_pred)
```

```
print('The accuracy of prediction is',accscore*100)
```

OUTPUT

PassiveAggressiveClassifier

```
[ ] from sklearn.linear_model import PassiveAggressiveClassifier  
  
    model = PassiveAggressiveClassifier(max_iter=50)  
    model.fit(tfidf_train, y_train)
```

```
PassiveAggressiveClassifier(max_iter=50)
```

```
[▶] y_pred = model.predict(tfidf_test)  
     accscore = accuracy_score(y_test, y_pred)  
     print('The accuracy of pediction is',accscore*100)
```

```
[▶] The accuracy of pediction is 99.59020044543429
```


CONCLUSION

In conclusion, a model to discriminate between fake and real news articles has been created using natural language processing (NLP) approaches.

The project involved pre-processing data by tokenizing, stemming, and removing stop words. The Tfidf Vectorizer was used to vectorize the dataset once it had been divided into training and testing sets. On the basis of their accuracy score, two machine learning models—Logistic Regression and Passive Aggressive Classifier—were trained and assessed. The model's high accuracy rating suggests that NLP methods can successfully identify bogus news. By adding more sophisticated NLP methods and expanding the dataset, further advancements can be realised. In general, this project shows how useful NLP can be for journalism and information integrity.