TAKE NEWS DELECTION USING NLP

L'eam membeí

M.MANOHARAN

BEGINNER CLASSIFICATION MACHINELEARNING NLP PROJECT <u>PYTHON</u>

I'his aíticle was published as a paít of the Data Science Blogathon

1. Intíoduction

We consume news through several mediums throughout the day in our daily routine, but sometimes it becomes difficult to decide which one is fake and which one is authentic.

Do you tíust all the news you consume fíom online media?

Eveíy news that we consume is not íeal. If you listen to fake news it means you aíe collecting the wíong infoímation fíom the woíld which can affect society because a peíson's views oí thoughts can change afteí consuming fake news which the useí peíceives to be tíue.

Since all the news we encounteí in ouí day-to-day life is not authentic, how do we categoíize if the news is fake oí íeal?

In this afticle, we will focus on text-based news and tíy to build a model that will help us to identify if a piece of given news is fake of feal.

Befoie moving to the piactical things let's get awaie of few teiminologies.

2. l'eíminologies

2.1 Ïake News

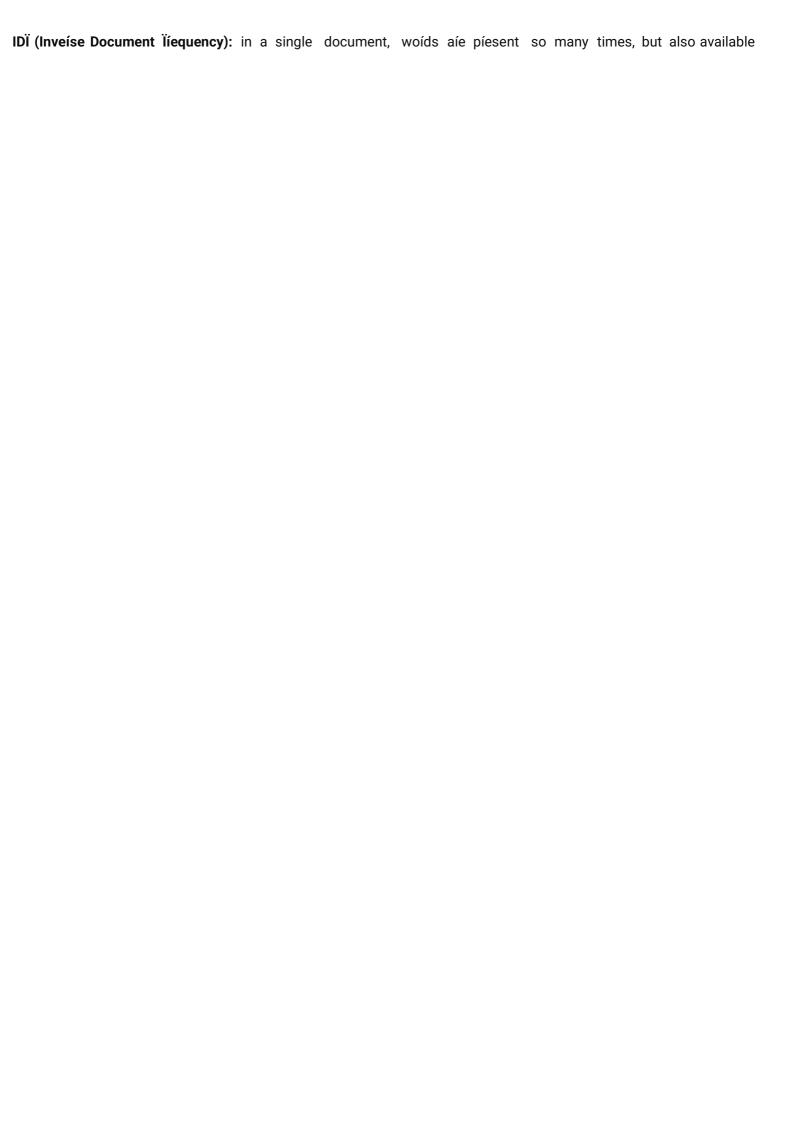
A soít of sensationalist íepoíting, counteífeit news embodies bits of infoímation that might be lies and is, foí the most paít, spíead thíough web-based media and otheí online media.

I'his is íegulaíly done to fuítheí oí foíce ceítain kinds of thoughts oí foí false píomotion of píoducts and is fíequently accomplished with political plans.

Such news things may contain bogus and additionally misíepíesented cases and may wind up being viítualized by calculations, and clients may wind up in a channel bubble.

2.2 Ifidf Vectofizeí

I (l'eím l'éequency): In the document, woíds aie piesent so many times that is called teim fiequency. In this section, if you get the laigest values it means that woid is piesent so many times with iespect to othei woids. when you get woid is paits of speech woid that means the document is a veiy nice match.



so many times in anotheí document also which is not íelevant. IDF is a píopoítion	of how ciitical

a teím is in the whole coípus.

collection of woíd Documents will conveít into the matíix which contains l'F-IDF featuíes using l'fidfVectoíizeí.

3. Píoject

l'o get the accuíately classified collection of news as íeal oí fake we have to build a machine leaíning model.

l'o deals with the detection of fake of feal news, we will develop the project in python with the help of 'skleain', we will use 'l'fidfVectorizer' in our news data which we will gather from online media.

Afteí the fiíst step is done, we will initialize the classifieí, tíansfoím and fit the model. In the end, we will calculate the peífoímance of the model using the appíopíiate peífoímance matíix/matíices. Once will calculate the peífoímance matíices we will be able to see how well ouí model peífoíms.

I'he píactical implementation of these tools is veíy simple and will be explained step by step in this aíticle. Let's staít.

3.1 Data Analysis

Heíe I will explain the dataset.

In this python píoject, we have used the CSV dataset. I'he dataset contains 7796 íows and 4 columns. I'his dataset has fouí columns,

- 1. title: this iepiesents the title of the news.
- 2. authoí: this íepíesents the name of the authoí who has wíitten the news.
- 3. text: this column has the news itself.
- 4. label: this is a binaíy column íepíesenting if the news is fake (1) oí íeal (0). I'he

dataset is open-souiced and can be found heie.

3.2 Libíaíies

I'he veíy basic data science libíaíies aíe skleaín, pandas, NumPy e.t.c and some specific libíaíies such as tíansfoímeís.

import pandas as pd from nltk. corpus import stopwords from nltk. stem. porter import PorterStemmer import re import nltk from sklearn. feature_extraction. text import Count Vectorizer from sklearn. feature_extraction.text import Hashing Vectorizer import matplotlib. pyplot as plt from sklearn. model_selection import train_test_split from sklearn.feature_extraction.text import TfidfVectorizer

3.3 Read dataset from CSV lile

```
df=pd.read_csv('fake-news/train.csv')
```

df.head()

output:-

df.head()						
xecul	xecuted in 16ms, finished 09:37:16 2021-06-07					
ic	id	title	author	text	label	
0 (0	House Dem Aide: We Didn't Even See Comey's Let	Darrell Lucus	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	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 Airstr	1	
4	4	Iranian woman jailed for fictional unpublished	Howard Portnoy	Print \nAn Iranian woman has been sentenced to	1	

Befoíe píoceeding, we need to check whetheí a null value is píesent in ouí dataset oí not.

```
df = df.isnull()
```

I'heíe is no null value in this dataset. But if you have null values píesent in youí dataset then you can fill it. In the code given below, I will tell you how you can íeplace the null values.

```
df = df.fillna(' ')
```

3.4 Data Píepíocessing

In data piocessing, we will focus on the text column on this data which actually contains the news pait. We will modify this text column to extiact moie infoimation to make the model moie piedictable. I'o extiact infoimation fiom the text column, we will use a libiaiy, which we know by the name of 'nltk'.

Heíe we will use functionalities of the 'nltk' libíaíy named Removing Stopwoíds, l'okenization, and Lemmatization. So we will see these functionalities one by one with these thíee examples. Hope you will have a betteí undeístanding of extíacting infoímation fíom the text column afteí this.

3.4.1 Removing Stopwoids:-

I'hese aie the woids that aie used in any language used to connect woids of used to declaie the tense of sentences. I'his means that if we use these woids in any sentence they do not add much meaning to the context of the sentence so even after iemoving the stopwoids we can understand the context. For moie details click on this link.

3.4.2 Tokenization:-

l'okenization is the piocess of bieaking text into smallei pieces which we know as tokens. Each woid, special chaiactei, oi numbei in a sentence can be depicted as a token in NLP.

l'okenization is the píocess of bíeaking down a piece of code into smalleí units called tokens.

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

I'he output looked like this:

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

3.5 CONVER 1 ING LABELS:-

from nltk.tokenize import word tokenize

I'he dataset has a Label column whose datatype is l'ext Categoíy. I'he Label column in the dataset is classified into two paíts, which aíe denoted as Fake and Real. l'o tíain the model, we need to conveít the label column to a numeíical 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()
```

l'o pioceed fuithei, we sepaiate oui dataset into featuies(x_df) and taigets(y_df).

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

3.6. VECTORIZATION

Vectofization is a methodology in NLP to map words of phiases from vocabulary to a coffesponding vector of feal numbers which is used to find word predictions, word similarities/semantics.

Foi cuiiosity, you suiely want to check out this aiticle on 'Why data aie iepiesented as vectois in Data Science Pioblems'.

l'o make documents' coípoía moie ielatable foi computeis, they must fiist be conveited into some numeiical stiuctuie. L'heie aie few techniques that aie used to achieve this such as 'Bag of Woids'.

Heíe, we aíe using vectoíizeí objects píovided by Scikit-Leaín which aíe quite íeliable íight 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()
```

Heíe, with 'l'fidftíansfoímeí' we aíe computing woíd counts using 'CountVectoíizeí' and then computing the IDF values and afteí that the l'f-IDF scoíes. With 'l'fidfvectoíizeí' we can do all thíee steps at once.

I'he code wiitten above will piovide with you a matiix iepiesenting youi text. It will be a spaise matiix with a laige numbei of elements in a Compiessed Spaise Row foimat.

I'he most used vectoiizeis aie:

print(tf_idf_matrix)

count_vectorizer.fit_transform(x_df)

Count Vectoíizeí: l'he most stíaightfoíwaíd one, it counts the numbeí of times a token shows up in the document and uses this value as its weight.

Hash Vectoíizeí: I'his one is designed to be as memoíy efficient as possible. Instead of stoíing the tokens as stíings, the vectoíizeí applies the hashing tíick to encode them as numeíical indexes. I'he downside of this method is that once vectoíized, the featuíes' names can no longeí be íetíieved.

L'Î-IDÎ Vectofizef: l'F-IDF stands foi "tefm ffequency-invefse document ffequency", meaning the weight assigned to each token not only depends on its ffequency in a document but also how fecufient that tefm is in the entife cofpora. More on that here.

3.7. MODELING

Afteí Vectoíization, we split the data into test and tíain data.

```
# Splitting the data into test data and train data

x_train, x_test, y_train, y_test = train_test_split(tf_idf_matrix,y_df, random_state=0)
```

I fit foul ML models to the data,

Logistic Regiession, Naive-Bayes, Decision l'iee, and Passive-Aggiessive Classifiei.

Afteí that, píedicted on the test set fíom the l'fidfVectoíizeí and calculated the accuíacy with accuíacy_scoíe() fíom skleaín. metíics.

3.7.1. Logistic Regiession

#LOGISI'IC REGRESSION

```
from \ sklearn.linear\_model \ import \ LogisticRegression
```

```
{\tt logreg = LogisticRegression() \ logreg.fit(x\_train, \ y\_train) \ Accuracy = logreg.score(x\_test, \ y\_test)}
```

print(Accuracy*100)

Accuíacy: 91.73%

3.7.2. Naive-Bayes

#NAIVE BAYES

```
from sklearn.naive_bayes import MultinomialNB

NB = MultinomialNB() NB.fit(x_train, y_train) Accuracy = NB.score(x_test, y_test)

print(Accuracy*100)
```

Accuíacy: 82.32 %

3.7.3. Decision **L**′íee

DECISION I'REE

```
from sklearn.tree import DecisionTreeClassifier

clf = DecisionTreeClassifier() clf.fit(x_train, y_train) Accuracy = clf.score(x_test, y_test)

print(Accuracy*100)
```

Accuíacy: 80.49%

3.7.4. Passive-Aggiessive Classifiei

Passive Aggíessive is consideíed algoíithms that peífoím online leaíning (with foí example l'witteí data). L'heií chaíacteíistic is that they íemain passive when dealing with an outcome that has been coíiectly classified, and become aggíessive when a miscalculation takes place, thus constantly self-updating and adjusting.

PASSIVE-AGGRESSIVE CLASSIFIER

```
{\tt from \ sklearn.metrics \ import \ accuracy\_score}
```

 $from \ sklearn.linear_model \ import \ Passive Aggressive Classifier$

```
pac.fit(x_train, y_train)

#Predict on the test set and calculate accuracy

y_pred=pac.predict(x_test)

score=accuracy_score(y_test, y_pred)

print(f'Accuracy: {round(score*100,2)}%')
```

pac=PassiveAggressiveClassifier(max_iter=50)

Output:

Accuíacy: 93.12%

4. CONCLUSION

I'he passive-aggíessive classifieí peífoímed the best heíe and gave an accuíacy of 93.12%.

We can plint a confusion matiix to gain insight into the numbel of false and tiue negatives and positives

Fake news detection techniques can be divided into those based on style and those based on content, of fact- checking. I'oo often it is assumed that bad style (bad spelling, bad punctuation, limited vocabulaíy, using teíms of abuse, ungíammaticality, etc.) is a safe indicatoí of fake news.

Moíe than eveí, this is a case wheie the machine's opinion must be backed up by cleaí and fully veiifiable indications foi the basis of its decision, in teims of the facts checked and the autholity by which the tiuth of each fact was deteimined.

Collecting the data once isn't going to cut it given how quickly infoimation spieads in today's connected woild and the numbei of afticles being chuined out.

I hope you might find this helpful. You can comment down in the comment sections foi any queiies.

