DATE:18-10-2023

PHASE-3:DEVELOPMENT

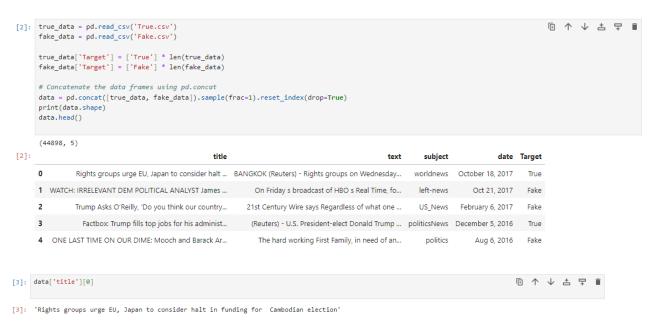
AIM:

Load the fake news dataset and preprocess the textual data.

IMPORTING LIBRARIES

```
[1]: # import modules
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import re
import nltk
import varnings
import sklearn
%matplotlib inline
warnings.filterwarnings('ignore')
```

LOADING THE DATASET



```
[4]: data['text'][0]
```

[4]: 'data['text'][0]

[4]: 'BANGKOK (Reuters) - Rights groups on Wednesday urged the European Union and Japan to consider halting their funding for the election panel in Cambodia, if the ruling party succeeds in a bid to dissolve the main opposition party ahead of next year s general election. The ruling Cambodia People s Party (C PP) has launched a crackdown on its cristics, including politicians, indepentent media and non-government bodies. Nearly half the opposition members of p arliament have fled abroad since September. In a session boycotted by the opposition, Cambodia s parliament voted on Monday to change party laws to red istribute seats if a party is dissolved. The measure came after the government filed a lawsuit this month seeking to dissolve the main opposition Cambodi an National Rescue Party (SURP). If the government s position to dissolve the opposition Cambodi national Rescue Party succeeds, next year s election will be a joke, Phil Robertson, deputy director for Asia at New York-based group Human Rights Watch, told Reuters. At that point, both the EU and Japan is should face reality and terminate their financial and technical assistance to avoid lending credibility to what will be a charade of democracy, he ad ded, speaking after a news conference in Bangkok. Japan and the EU are the two biggest foreign funders of the 2018 vote. China and the United States have also contributed, with the United States providing trucks and technical support, while Japan has given computers. Japan s embassy in Phnom Penh did no treply to a Reuters request for comment on the matter. George Edgar, head of the EU delegation to Cambodia, said the EU remains ready to support a credible electoral process but added that the polls should only go shead with the opposition sinvolvement. He urged Cambodian authorities not to go ahead with the dissolution of the opposition party. The EU remains ready to support a credible electoral process up to the National Assembly election in 2018. However we do not believe that a proce

```
[5]: data.info()
      <class 'pandas.core.frame.DataFrame'>
      RangeIndex: 44898 entries, 0 to 44897
Data columns (total 5 columns):
       # Column Non-Null Count Dtype
       0 title 44898 non-null object
           text
                      44898 non-null object
            subject 44898 non-null object
       3 date 44898 non-null object
4 Target 44898 non-null object
      dtypes: object(5)
          nory usage: 1.7+ MB
```

PREPROCESSING:

REMOVING NULL VALUES:

(from the information above there is no null values so nothing is removed)

```
[6]: #preprocessing
      #dron nul.l. val.ues
      data=data.dropna(axis=0)
[7]: len(data)
[7]: 44898
```

CONVERTING ALL STRINGS TO LOWERCASE:

```
[8]: #converting all strings to lowercase
      data['clean_news']=data['text'].str.lower()
      data['clean news']
[8]: 0
                 bangkok (reuters) - rights groups on wednesday...
                 on friday s broadcast of hbo s real time, fo...
21st century wire says regardless of what one ...
                 (reuters) - u.s. president-elect donald trump ...
                 the hard working first family, in need of an..
      44893
               21st century wire says does the american ideal...
                 barinas, venezuela (reuters)
                 phnom penh (reuters) - cambodian prime ministe...
                 geneva (reuters) - the united states wants to ...
beijing (reuters) - u.s. president donald trum...
      Name: clean_news, Length: 44898, dtype: object
```

REMOVING SPECIAL CHARACTERS , EXTRA SPACES AND ESCAPE CHARACTERS

```
[9]: #removing special characters , extra spaces and escape characters
data['clean_news']=data['clean_news'].str.replace('[^A-Za-Za-9-9\s]','')
data['clean_news']=data['clean_news'].str.replace('[\n]','')
data['clean_news']=data['clean_news'].str.replace('[\s+]','')
data['clean_news']

[9]: 0 bangkok (reuters) - rights groups on wednesday...
1 on friday s broadcast of hbo s real time, fo...
2 21st century wire says regardless of what one ...
3 (reuters) - u.s. president-elect donald trump ...
4 the hard working first family, in need of an...

44893 21st century wire says does the american ideal...
44894 barinas, venezuela (reuters) - tirelessly trav...
44895 phnom penh (reuters) - cambodian prime ministe...
44896 geneva (reuters) - te united states wants to ...
44897 beijing (reuters) - u.s. president donald trum...
Name: clean_news, Length: 44898, dtype: object
```

REMOVING STOP WORDS

```
[10]: nltk.download('stopwords')
       [nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\rsriv\AppData\Roaming\nltk_data...
       [nltk_data] Package stopwords is already up-to-date!
[10]: True
       from nltk.corpus import stopwords
       stop=stopwords.words('english')
       data['clean_news']=data['clean_news'].apply(lambda x: " ".join([word for word in x.split() if word not in stop]))
                                                                                                                            date Target
                                                                                                        subject
       0 Rights groups urge EU, Japan to consider halt
                                                          BANGKOK (Reuters) - Rights groups on
                                                                                                                     October 18,
                                                                                                                                                  bangkok (reuters) - rights groups
                                                                                                    worldnews
                                                                                                                                    True
       1 WATCH: IRRELEVANT DEM POLITICAL ANALYST
                                                          On Friday s broadcast of HBO s Real Time,
                                                                                                                                              friday broadcast hbo real time, former
                                                                                                                    Oct 21, 2017 Fake
                                                          21st Century Wire says Regardless of what
                   Trump Asks O'Reilly, 'Do you think our
                                                                                                                                               21st century wire says regardless one
       2
                                                                                                      US News February 6, 2017 Fake
                                                            (Reuters) - U.S. President-elect Donald Trump ... politicsNews
                                                                                                                                               (reuters) - u.s. president-elect donald
       3 Factbox: Trump fills top jobs for his administ...
                                                                                                                           2016
              ONE LAST TIME ON OUR DIME: Mooch and People Ar. The hard working First Family, in need of an...
                                                                                                                                             hard working first family, need another
                                                                                                       politics
                                                                                                                    Aug 6, 2016 Fake
```

TOKENIZATION:

```
[13]: nltk.download('punkt')

[nltk_data] Downloading package punkt to
   [nltk_data] C:\Users\rsriv\AppData\Roaming\nltk_data...
   [nltk_data] Package punkt is already up-to-date!

[13]: True

[14]: #Tokenization
   from nltk.tokenize import word_tokenize
    data['tokenized_news'] = data['clean_news'].apply(lambda x: word_tokenize(x))
    data_bead()
```

14]:	title	text	subject	date	Target	clean_news	tokenized_news
0	Rights groups urge EU, Japan to consider halt	BANGKOK (Reuters) - Rights groups on Wednesday	worldnews	October 18, 2017	True	bangkok (reuters) - rights groups wednesday ur	[bangkok, (, reuters,), -, rights, groups, we
1	WATCH: IRRELEVANT DEM POLITICAL ANALYST James	On Friday s broadcast of HBO s Real Time, fo	left-news	Oct 21, 2017	Fake	friday broadcast hbo real time, former clinton	[friday, broadcast, hbo, real, time, ,, former
2	Trump Asks O'Reilly, 'Do you think our country	21st Century Wire says Regardless of what one	US_News	February 6, 2017	Fake	21st century wire says regardless one thinks d	[21st, century, wire, says, regardless, one, t
3	Factbox: Trump fills top jobs for his administ	(Reuters) - U.S. President-elect Donald Trump	politicsNews	December 5, 2016	True	(reuters) - u.s. president-elect donald trump	[(, reuters,), -, u.s., president- elect, dona
4	ONE LAST TIME ON OUR DIME: Mooch and Barack Ar	The hard working First Family, in need of an	politics	Aug 6, 2016	Fake	hard working first family, need another taxpay	[hard, working, first, family, ,, need, anothe

LEMMATIZATION:

]: n]	nltk.download('wordnet')										
[r	nltk_data] Downloading pac nltk_data] C:\Users\rs nltk_data] Package wordr										
]: Tr	rue										
fr le de	### ##################################										
]:	title	text	subject	date	Target	clean_news	tokenized_news	lemmatized_news			
]: 0	Rights groups urge EU, Japan to consider halt	text BANGKOK (Reuters) - Rights groups on Wednesday	subject worldnews	October 18, 2017	Target True	clean_news bangkok (reuters) - rights groups wednesday ur	tokenized_news [bangkok, (, reuters,), -, rights, groups, we	[bangkok, (, reuters,), -, right, group, wedn			
_	Rights groups urge EU,	BANGKOK (Reuters) - Rights groups on		October 18,		bangkok (reuters) - rights groups	[bangkok, (, reuters,), -,	[bangkok, (, reuters,), -, right, group, wedn			
_	Rights groups urge EU, Japan to consider halt WATCH: IRRELEVANT DEM POLITICAL ANALYST James Trump Adds Ciffeills . Des	BANGKOK (Reuters) - Rights groups on Wednesday On Friday's broadcast of	worldnews	October 18, 2017	True	bangkok (reuters) - rights groups wednesday ur friday broadcast hbo real	[bangkok, (, reuters,), -, rights, groups, we	[bangkok, (, reuters,), -, right, group, wedn			
0	Rights groups urge EU, Japan to consider halt WATCH: IRRELEVANT DEM POLITICAL ANALYST James Trump Asks O'Reilly, 'Do	BANGKOK (Reuters) - Rights groups on Wednesday On Friday's broadcast of HBO's Real Time, fo 21st Century Wire says	worldnews left-news	October 18, 2017 Oct 21, 2017 February 6,	True Fake	bangkok (reuters) - rights groups wednesday ur friday broadcast hbo real time, former clinton 21st century wire says	[bangkok, (, reuters,), -, rights, groups, we [friday, broadcast, hbo, real, time, ,, former [21st, century, wire, says,	[bangkok, (, reuters,), -, right, group, wedn [friday, broadcast, hbo, real, time, ,, former [21st, century, wire, say,			

CREATE SENTENCES TO GET CLEAN TEXT AS INPUT FOR VECTORS

	<pre>data['clean_text'] = data['lemmatized_news'].apply(lambda x : return_sentences(x)) data.head()</pre>									
	title	text	subject	date	Target	clean_news	tokenized_news	lemmatized_news	clean_text	
0	Rights groups urge EU, Japan to consider halt	BANGKOK (Reuters) - Rights groups on Wednesday	worldnews	October 18, 2017	True	bangkok (reuters) - rights groups wednesday ur	[bangkok, (, reuters,), -, rights, groups, we	[bangkok, (, reuters,), -, right, group, wedn	bangkok (reuters - right group wednesday ur	
1	WATCH: IRRELEVANT DEM POLITICAL ANALYST James	On Friday s broadcast of HBO s Real Time, fo	left-news	Oct 21, 2017	Fake	friday broadcast hbo real time, former clinton	[friday, broadcast, hbo, real, time, ,, former	[friday, broadcast, hbo, real, time, ,, former	friday broadcas hbo real time former clinto.	
2	Trump Asks O'Reilly, 'Do you think our country	21st Century Wire says Regardless of what one	US_News	February 6, 2017	Fake	21st century wire says regardless one thinks d	[21st, century, wire, says, regardless, one, t	[21st, century, wire, say, regardless, one, th	21st century wir say regardless on think don.	
3	Factbox: Trump fills top jobs for his administ	(Reuters) - U.S. President-elect Donald Trump	politicsNews	December 5, 2016	True	(reuters) - u.s. president-elect donald trump	[(, reuters,), -, u.s., president-elect, dona	[(, reuters,), -, u.s., president-elect, dona	(reuters) - u.s president-elec donald trum.	
4	ONE LAST TIME ON OUR DIME: Mooch and Barack Ar	The hard working First Family, in need of an	politics	Aug 6, 2016	Fake	hard working first family, need another taxpay	[hard, working, first, family, ,, need, anothe	[hard, working, first, family, ,, need, anothe	hard working firs family , need another taxpa.	

PREPARE DATA FOR THE MODEL. CONVERT LABEL IN TO BINARY

	<pre>data['Target'] = [1 if x == 'Fake' else 0 for x in data['Target']] data.head()</pre>											
	title text subject date Target clean_news tokenized_news lemmatized_news cle											
0	Rights groups urge EU, Japan to consider halt	BANGKOK (Reuters) - Rights groups on Wednesday	worldnews	October 18, 2017	0	bangkok (reuters) - rights groups wednesday ur	[bangkok, (, reuters,), -, rights, groups, we	[bangkok, (, reuters,), -, right, group, wedn	bangkok (reuters - right grou wednesday ur.			
1	WATCH: IRRELEVANT DEM POLITICAL ANALYST James	On Friday s broadcast of HBO s Real Time, fo	left-news	Oct 21, 2017	1	friday broadcast hbo real time, former clinton	[friday, broadcast, hbo, real, time, ,, former	[friday, broadcast, hbo, real, time, ,, former	friday broadcas hbo real time former clinto.			
2	Trump Asks O'Reilly, 'Do you think our country	21st Century Wire says Regardless of what one	US_News	February 6, 2017	1	21st century wire says regardless one thinks d	[21st, century, wire, says, regardless, one, t	[21st, century, wire, say, regardless, one, th	21st century wir say regardless on think don.			
3	Factbox: Trump fills top jobs for his administ	(Reuters) - U.S. President-elect Donald Trump	politicsNews	December 5, 2016	0	(reuters) - u.s. president-elect donald trump	[(, reuters,), -, u.s., president-elect, dona	[(, reuters,), -, u.s., president-elect, dona	(reuters) - u.s president-elec donald trum.			
4	ONE LAST TIME ON OUR DIME: Mooch and Barack Ar	The hard working First Family, in need of an	politics	Aug 6, 2016	1	hard working first family, need another taxpay	[hard, working, first, family, ,, need, anothe	[hard, working, first, family, ,, need, anothe	hard working firs family , nee another taxpa			

SPLIT THE DATASET

```
[22]:

X_train, X_test, y_train, y_test = train_test_split(data['clean_text'], data['Target'], test_size=0.2, random_state=5)

print(X_train.shape)
print(X_test.shape)

(35918,)
(8980,)
```

FEATURE EXTRACTION

(words to vectors)

Count vectorizer which considers the frequency of occurrence of a word across the corpus.

```
[23]: from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

# Assuming 'lemmatized_news' is a column in your DataFrame 'data'
count_vectorizer = CountVectorizer()
X = count_vectorizer.fit_transform(data['clean_text'])

# Get feature names
feature_names_count = count_vectorizer.get_feature_names_out()

print("CountVectorizer feature names:", feature_names_count)

CountVectorizer feature names: ['00' '000' '000' '... 'zzzzzzzzzzzzz' 'émigré']
```

TF-IDF: Term Frequency - Inverse Document Frequency

The term frequency is the number of times a term occurs in a document. Inverse document frequency is an inverse function of the number of documents in which that a given word occurs.

The product of these two terms gives tf-idf weight for a word in the corpus.

CONCLUSION:

In the preprocessing stage of the fake news detection dataset, a series of crucial steps were performed using NLP techniques. Punctuation and stopwords were removed to clean the text data, and tokenization was applied, ensuring that each document was divided into individual words. Further refinement was achieved through lemmatization, standardizing words to their base form. Feature vectors were then generated, capturing the essence of the text data. Additionally, it is important to highlight that using sparse representations for text data, like TF-IDF, can significantly reduce memory usage and improve model performance.