```python
 1   import pandas as pd
 2  import numpy as np
 3  import matplotlib.pyplot as plt
 4  import seaborn as sns
 5  data = pd.read_csv("fake_news.csv")
 6  data.head()
 7  data.shape
 8  data.info()
 9  data.isna().sum()
10  data = data.drop(['id'] , axis=1)
11   data = data.fillna('')
12  data['content'] = data['author'] +''+ data['title']+''+data['text']
13   data = data.drop(['title','author','text'], axis=1)
14  data.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20800 entries, 0 to 20799
Data columns (total 5 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   id      20800 non-null  int64
 1   title   20242 non-null  object
 2   author  18843 non-null  object
 3   text    20761 non-null  object
 4   label   20800 non-null  int64
dtypes: int64(2), object(3)
memory usage: 812.6+ KB
```

Out[3]:

| | label | content |
|---|---|---|
| **0** | 1 | Darrell LucusHouse Dem Aide: We Didn't Even Se... |
| **1** | 0 | Daniel J. FlynnFLYNN: Hillary Clinton, Big Wom... |
| **2** | 1 | Consortiumnews.comWhy the Truth Might Get You ... |
| **3** | 1 | Jessica Purkiss15 Civilians Killed In Single U... |
| **4** | 1 | Howard PortnoyIranian woman jailed for fiction... |

```
In [4]:  1  data['content'] = data['content'].apply(lambda x: " ".join(x.lower() for x in x.split()))
```

```
In [5]:  1  data['content'] = data['content'].str.replace('[^\w\s]','')
```

C:\Users\Sai Krishna Hari\AppData\Local\Temp\ipykernel_5984\3643324700.py:1: FutureWarning: The default valu
e of regex will change from True to False in a future version.
  data['content'] = data['content'].str.replace('[^\w\s]','')

```
In [6]:  1  import nltk
         2  nltk.download("stopwords")
```

[nltk_data] Downloading package stopwords to C:\Users\Sai Krishna
[nltk_data]     Hari\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!

Out[6]:  True

```
In [8]:  1  from nltk.corpus import stopwords
         2  stop=stopwords.words('english')
         3  data['content']=data['content'].apply(lambda x:" ".join(x for x in x.split() if x not in stop))
```

```
In [11]:  1  from nltk.stem import WordNetLemmatizer
          2  from textblob import Word
          3  data['content']=data['content'].apply(lambda x:"".join([Word(word).lemmatize() for word in x.split()]))
          4  data['content'].head()
```

Out[11]:  0    darrelllucushousedemaidedidntevenseecomeyslett...
          1    danieljflynnflynnhillaryclintonbigwomancampusb...
          2    consortiumnewscomwhytruthmightgetfiredwhytruth...
          3    jessicapurkiss15civiliankilledsingleuairstrike...
          4    howardportnoyiranianwomanjailedfictionalunpubl...
          Name: content, dtype: object

```
In [12]:  1  x=data[['content']]
          2  y=data['label']
```

```python
In [13]:  1  from sklearn.model_selection import train_test_split
```

```python
In [15]:  1  x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=45,stratify=y)
```

```python
In [16]:  1  print(x_train.shape)
```

(20755, 1)

```python
In [17]:  1  print(x_test.shape)
```

(45, 1)

```python
In [18]:  1  print(y_train.shape)
```

(20755,)

```python
In [19]:  1  print(y_test.shape)
```

(45,)

```python
In [21]:  1  from sklearn.feature_extraction.text import TfidfVectorizer
          2  tfidf_vect = TfidfVectorizer(analyzer='word', token_pattern=r'\w{1,}', max_features=5000)
          3  tfidf_vect.fit(data['content'])
          4  xtrain_tfidf = tfidf_vect.transform(x_train['content'])
          5  xtest_tfidf = tfidf_vect.transform(x_test['content'])
          6
```

```python
In [ ]:  1
```