

FAKE NEWS DETECTION USING NLP

**TEAM MEMBER:822621104307 –
P.THOOYAVAN**

Phase-2 Document submission



In this phase, we embark on an exciting journey to enhance the accuracy and robustness of our prediction system by delving into cutting-edge techniques such as ensemble methods and advanced deep learning architectures. Our primary objective is to elevate the efficiency of fake news detection, and to do so, we're considering the utilization of advanced models like Long Short-Term Memory (LSTM) networks and Bidirectional Encoder Representations from Transformers (BERT).

1. **Ensemble Methods:** Ensemble methods are a powerful way to combine the predictions of multiple models to improve overall accuracy. We can explore techniques such as bagging and boosting, where we train multiple fake news detection models and combine their outputs to make more informed predictions. This can help mitigate the weaknesses of individual models and enhance overall system performance.

2. **Deep Learning Models:**

- ****LSTM (Long Short-Term Memory)**:** LSTM networks are well-suited for sequential data, making them an excellent choice for text-based tasks like fake news detection. We can employ LSTM networks to capture the temporal dependencies in news articles, helping us identify patterns that might be indicative of fake news.

- ****BERT (Bidirectional Encoder Representations from Transformers)**:** BERT has revolutionized natural language processing tasks. Its bidirectional context understanding can significantly enhance our ability to discern the subtleties of language and context in news articles. Fine-tuning BERT for fake news detection can lead to a substantial improvement in accuracy.

3. **Data Augmentation:**

To further enrich our training data, we can employ data augmentation techniques. These methods involve generating synthetic examples based on our existing data, helping the models learn more diverse patterns and making them more robust against adversarial attacks.

4. **Transfer Learning:**

We can investigate the use of transfer learning from pre-trained models. By starting with a model trained on a large corpus of text data, we can save time and resources and fine-tune it for fake news detection. This can be especially beneficial if labeled fake news data is limited.

5. ****Regularization and Hyperparameter Tuning****:

Ensuring our models are well-regularized and fine-tuning hyperparameters is crucial for their performance. Techniques like dropout and batch normalization can help prevent overfitting, while grid search or Bayesian optimization can help find the best hyperparameter configurations.

6. ****Explainable AI (XAI)****:

It's essential to understand why a model makes a particular prediction, especially in sensitive applications like fake news detection. We can explore techniques for making deep learning models more interpretable, such as attention mechanisms or gradient-based attribution methods.

7. ****Continuous Monitoring and Feedback Loop****:

Once our enhanced models are deployed, it's crucial to continuously monitor their performance, collect feedback, and retrain them with fresh data. Fake news is a constantly evolving issue, and our models need to adapt to new trends and tactics used by malicious actors.

By venturing into these advanced techniques, we aim to fortify our fake news detection system, making it more resilient, accurate, and adaptable to the ever-changing landscape of misinformation. This proactive approach positions us at the forefront of the battle against fake news, safeguarding the integrity of information dissemination in the digital age.

SOURCE CODE:

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import plotly.express as px
import plotly.graph_objs as go
from plotly.subplots import make_subplots
```

```

import nltk
from nltk.corpus import stopwords
import tensorflow as tf
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import ModelCheckpoint
from sklearn.model_selection import train_test_split
from transformers import AutoTokenizer, TFAutoModelForSequenceClassification

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
nltk.download('stopwords')
/opt/conda/lib/python3.10/site-packages/scipy/__init__.py:146: UserWarning: A
NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (det
ected version 1.23.5
  warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}")
/kaggle/input/fake-and-real-news-dataset/True.csv
/kaggle/input/fake-and-real-news-dataset/Fake.csv
[nltk_data] Downloading package stopwords to /usr/share/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
Out[1]:
True
fake_news_path = "/kaggle/input/fake-and-real-news-dataset/Fake.csv"
real_news_path = "/kaggle/input/fake-and-real-news-dataset/True.csv"
fake_news = pd.read_csv(fake_news_path)
real_news = pd.read_csv(real_news_path)
linkcode
fake_news.head(3)

```

	title	text	subject	date
0	Donald Trump Sends Out Embarrassing New Year'...	Donald Trump just couldn t wish all Americans ...	News	December 31, 2017
1	Drunk Bragging Trump Staffer Started Russian ...	House Intelligence Committee Chairman Devin Nu...	News	December 31, 2017
2	Sheriff David Clarke Becomes An Internet Joke...	On Friday, it was revealed that former Milwauk...	News	December 30, 2017

```
real_news.head(3)
```

	title	text	subject	date
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people will...	politicsNews	December 29, 2017
2	Senior U.S. Republican senator: 'Let Mr. Muell...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017

```

real = real_news.copy()
fake = fake_news.copy()
real['Label'] = 'Real'
fake['Label'] = 'Fake'
linkcode
news = pd.concat([real, fake], axis=0, ignore_index=True)
news.reset_index()
news.head()

```

	title	text	subject	date	Label
0	As U.S. budget fight looms, Republicans flip t...	WASHINGTON (Reuters) - The head of a conservat...	politicsNews	December 31, 2017	Real
1	U.S. military to accept transgender recruits o...	WASHINGTON (Reuters) - Transgender people will...	politicsNews	December 29, 2017	Real
2	Senior U.S. Republican senator: 'Let Mr. Muell...	WASHINGTON (Reuters) - The special counsel inv...	politicsNews	December 31, 2017	Real
3	FBI Russia probe helped by Australian diplomat...	WASHINGTON (Reuters) - Trump campaign adviser ...	politicsNews	December 30, 2017	Real
4	Trump wants Postal Service to charge 'much mor...	SEATTLE/WASHINGTON (Reuters) - President Donal...	politicsNews	December 29, 2017	Real

```

print(f"Samples available: {news.shape[0]}\n#features of dataset: {news.shape[1]}")
Samples available: 44898
#features of dataset: 5
news_ds = news.sample(1000).drop(['title', 'date', 'subject'], axis=1)
news_ds.head(3)

```

	text	Label
36923	If local law enforcement begins to act like i...	Fake
38038	Obama and HUD want to give one last freebie to...	Fake

	text	Label
28271	Corey Lewandowski got some good news last week...	Fake

```

CLASS_NAMES = ['Fake', 'Real']
class_mapper = {
    'Fake':0,
    'Real':1
}
news_ds['Label'] = news_ds['Label'].map(class_mapper)
news_ds.head(3)

```

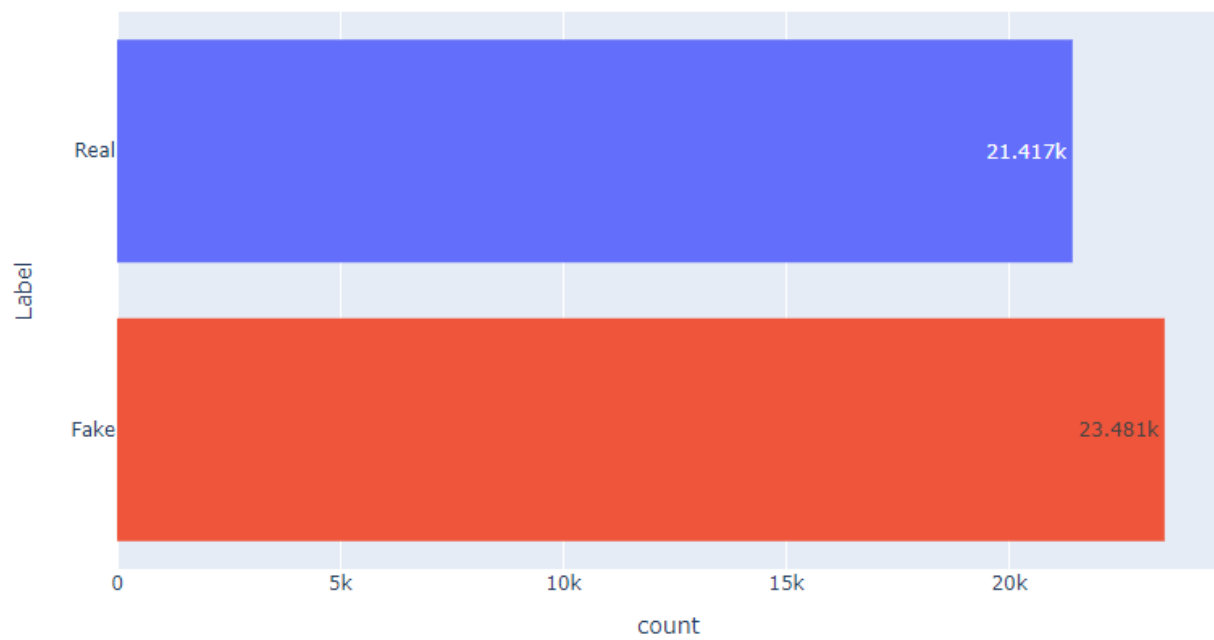
	text	Label
36923	If local law enforcement begins to act like i...	0
38038	Obama and HUD want to give one last freebie to...	0
28271	Corey Lewandowski got some good news last week...	0

```

class_dist = px.histogram(data_frame=news,
                           y='Label',
                           color='Label',
                           title='Fake vs Real news Original dataset',
                           text_auto=True)
class_dist.update_layout(showlegend=False)
class_dist.show()

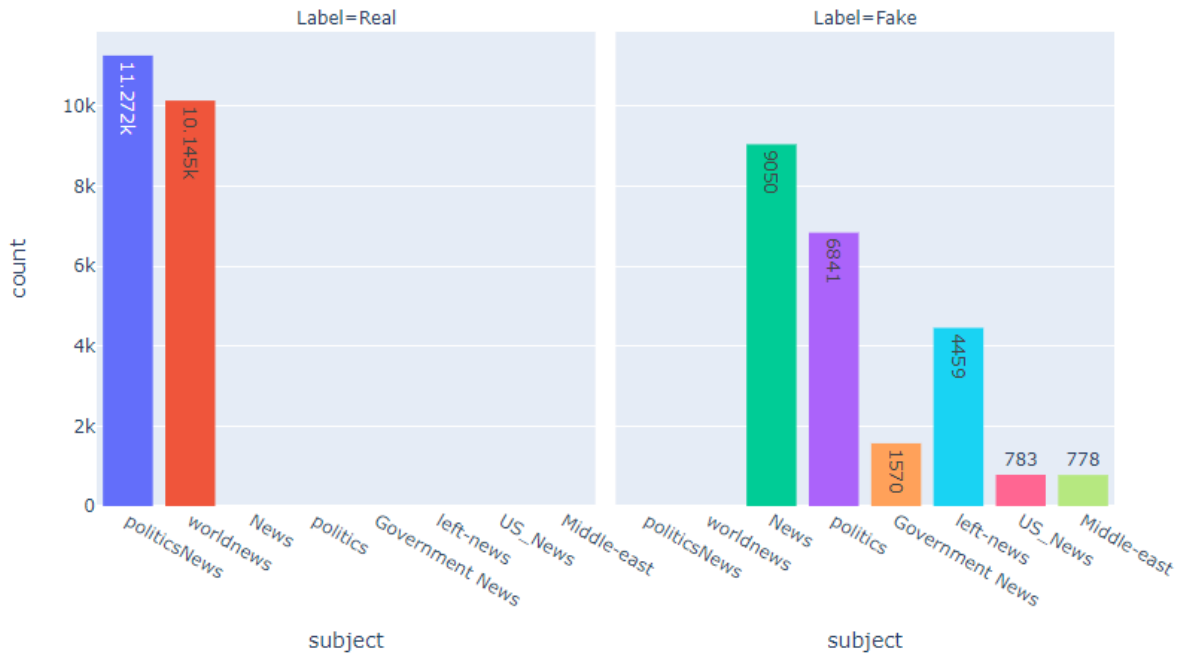
```

Fake vs Real news Original dataset

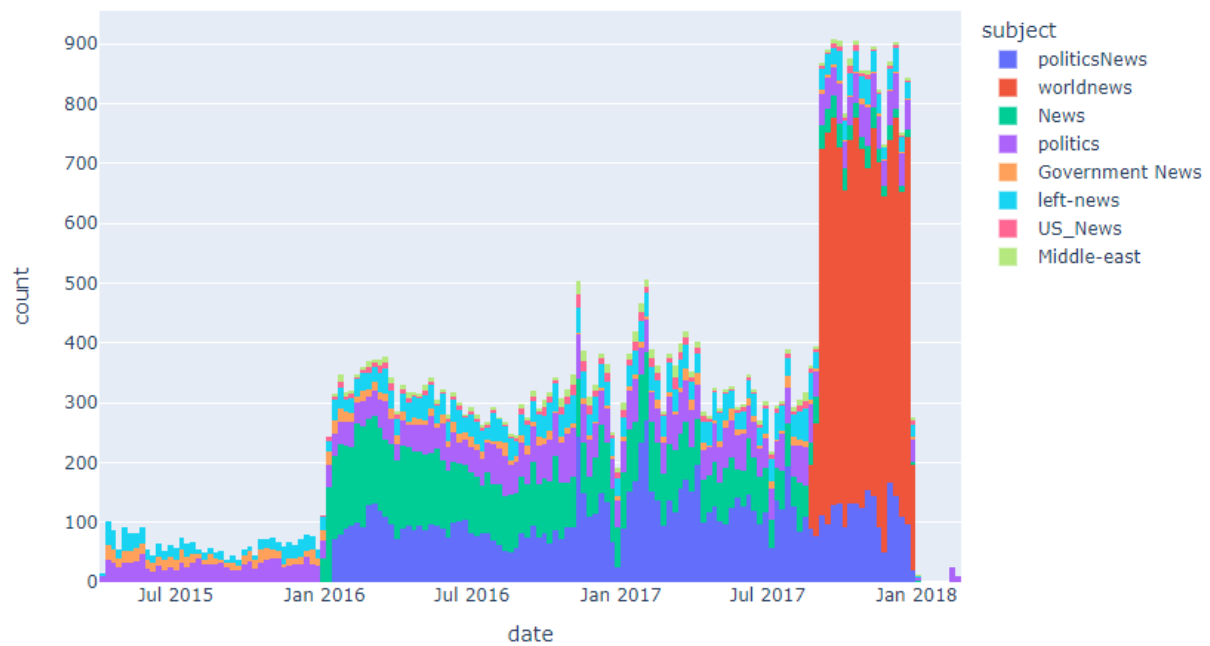


```
subject_dist = px.histogram(data_frame=news,  
                             x='subject',  
                             color='subject',  
                             title='Fake vs Real news Subject Distribution',  
                             text_auto=True,  
                             facet_col='Label')  
subject_dist.update_layout(showlegend=False)  
subject_dist.show()
```

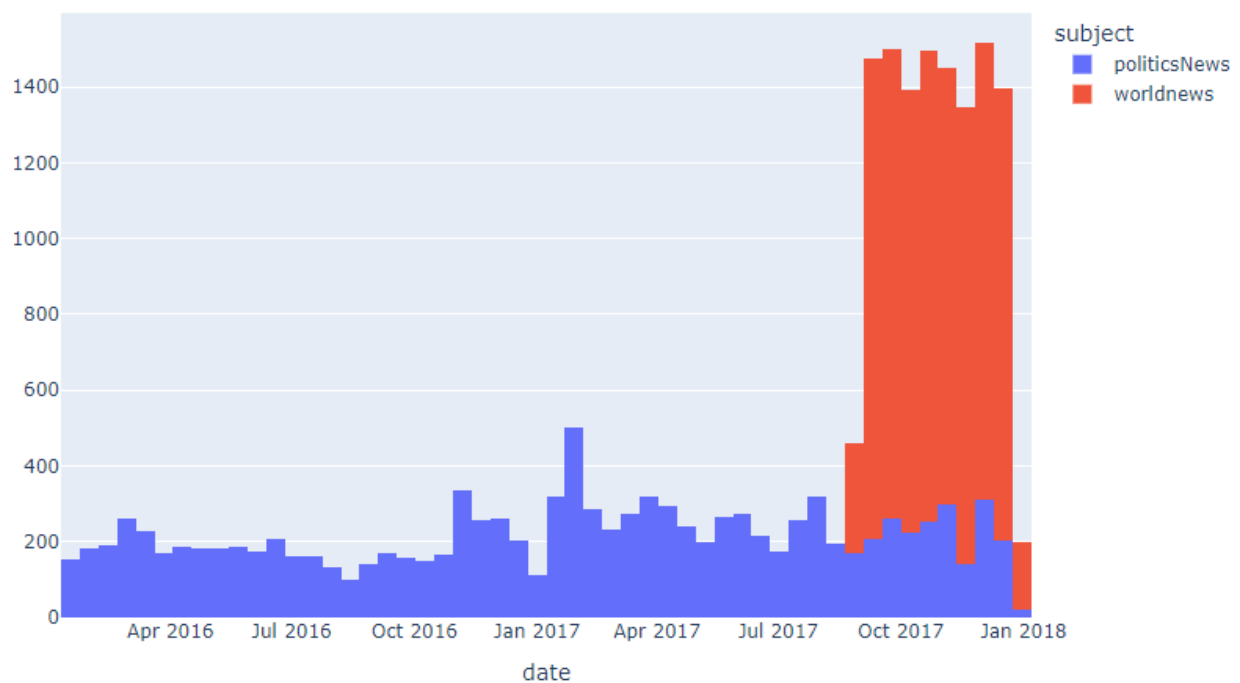
Fake vs Real news Subject Distribution



```
news.date.unique().max()
'https://fedup.wpengine.com/wp-content/uploads/2015/04/hillarystreetart.jpg'
list(filter(lambda x:len(x)>20, news.date.unique()))
['https://100percentfedup.com/served-roy-moore-vietnamletter-veteran-sets-
record-straight-honorable-decent-respectable-patriotic-commander-soldier/',
'https://100percentfedup.com/video-hillary-asked-about-trump-i-just-want-to-
eat-some-pie/',
'https://100percentfedup.com/12-yr-old-black-conservative-whose-video-to-oba
ma-went-viral-do-you-really-love-america-receives-death-threats-from-left/',
'https://fedup.wpengine.com/wp-content/uploads/2015/04/hillarystreetart.jpg'
,
'https://fedup.wpengine.com/wp-content/uploads/2015/04/entitled.jpg',
'MSNBC HOST Rudely Assumes Steel Worker Would Never Let His Son Follow in Hi
s Footsteps...He Couldn't Be More Wrong [Video]']
linkcode
news = news[news['date'].map(lambda x:len(x)) <= 20]
news.date = pd.to_datetime(news['date'], format='mixed')
news.head()
```

```
real_sub_dist = px.histogram(data_frame=news[news['Label']=='Real'],
                             x='date',
                             color='subject')
real_sub_dist.show()
```

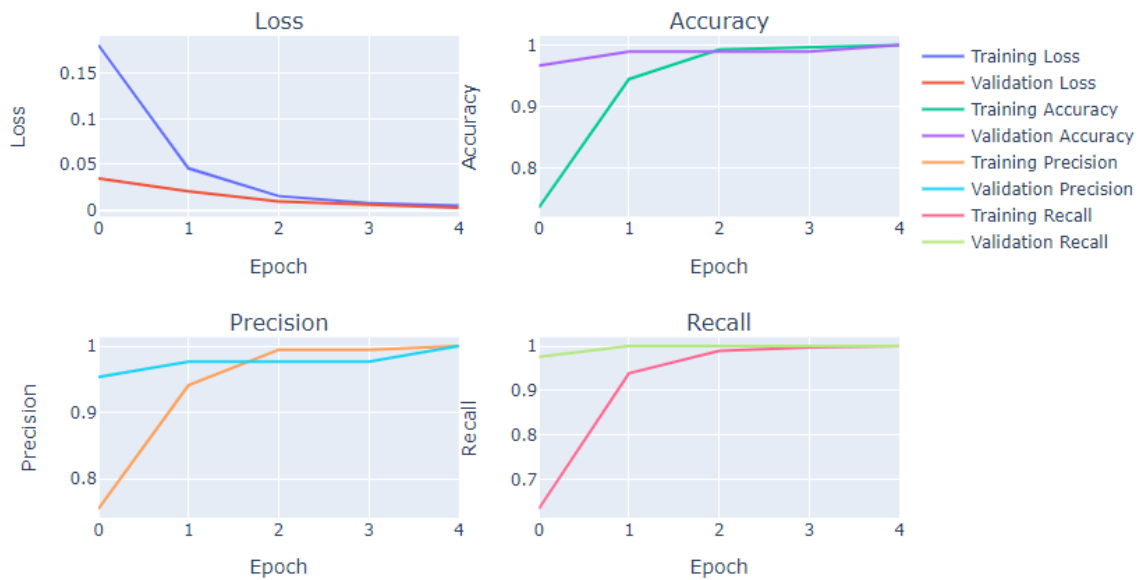


```
import string
```


model_history

	loss	Accuracy	Precision	Recall	val_loss	val_Accuracy	val_Precision	val_Recall
0	0.181528	0.735802	0.754777	0.633690	0.034657	0.966667	0.953488	0.97619
1	0.045669	0.944444	0.941019	0.938503	0.020463	0.988889	0.976744	1.00000
2	0.015262	0.992593	0.994624	0.989305	0.009228	0.988889	0.976744	1.00000
3	0.007200	0.996296	0.994667	0.997326	0.006398	0.988889	0.976744	1.00000
4	0.004542	1.000000	1.000000	1.000000	0.002476	1.000000	1.000000	1.00000

Model Training History



```
print(f"Test Loss      : {test_loss}")
print(f"Test Accuracy  : {test_acc}")
print(f"Test Precision   : {test_precision}")
print(f"Test Recall      : {test_recall}")
Test Loss      : 0.003277710871770978
Test Accuracy  : 1.0
Test Precision  : 1.0
Test Recall    : 1.0
def make_prediction(text, model=model):
    text = np.array([text])
    inputs = tokenize(text)
    return np.abs(np.round(model.predict(inputs, verbose=1).logits))

for _ in range(5):
    index = np.random.randint(test_X.shape[0])

    text = test_X[index]
    real = test_y[index]
    model_pred = make_prediction(text)
```

```
print(f"Original Text:\n\n{text}\n\nTrue: {CLASS_NAMES[int(real)]}\t\tPredicted: {CLASS_NAMES[int(model_pred)]}\n{'-'*100}\n")
```

CONCLUSION:

In the end, the application of NLP in fake news detection is not just a technological pursuit; it is a commitment to preserving the integrity of information, protecting individuals from the harmful effects of misinformation, and upholding the principles of truth and accuracy in communication. As NLP techniques continue to advance and adapt to the challenges of our digital age, we must remain steadfast in our commitment to using them responsibly and effectively for the betterment of society.