

Fake News Detection using Neural Network

Name: Sonu Sankhala

Reg. No. : 11805066

Course: INT246

CONTENTS

- Abstract
- Introduction
- Dataset Pre-Processing
- Model Construction & Discussion
- Conclusion
- References

Abstract

Fake news attracted attention both from the public and the academic communities and represents a phenomenon that has a significant impact on our social life, especially on the political world. Further more, fake news phenomenon provide an opportunity for malicious parties to manipulate public opinion and events such as elections. In this work, I propose a merged Machine learning model that detect fake articles regarding different characteristics. Therefore, I use some python technique and Passive Aggressive Classifiers algorithm to extract text based features and compare different architecture. I show on real dataset that the proposed approach is very efficient and allows to achieve high performances.

KEYWORD:- Fake news detection, Neural networks, Machine learning, Passive Aggressive Classifiers, Algorithm.

Introduction

Nowadays, fake news has become a common trend. Even trusted media houses are known to spread fake news and are losing their credibility. So, how can we trust any news to be real or fake? In the past decade, social media networks have become a valuable resource for data of different types, such as texts, photos, videos etc. On social media Fake news has experienced a resurgence of interest primarily due to the recent political climate and the growing concern around its negative effect. Even though the problem of fake news is not a new issue, detecting fake news is a complex task for us. In this project we deal with such fake news detection issue. In this project I've try to resolve that the given news fake or real, with the help of some algorithms using Neural Network.

Dataset Pre-Processing

Since fake news detection is relatively a new field, few datasets are available publicly. In my work, I conduct my experiments on the dataset proposed by kaggle.com . It contains about 20800 articles labeled for two classes : 10413 fake label and 10387 real label. The following attributes are included :

- id: unique id for a news article.
- title: the title of a news article.
- author: author of the news article.
- text: the text of the article; could be incomplete.
- label: a label that marks the article as potentially fake (1) or real (0).

In the pre-processing phase, we have removed from each article's text and title the punctuations and stop-word that are the most common words in a language like: "are", "as", "the" ..., ect.

```
In [2]: df = pd.read_csv('C:/Users/Sonu/fake_news_detector/train.csv')
conversion_dict = {0: 'Real', 1: 'Fake'}
df['label'] = df['label'].replace(conversion_dict)
df.label.value_counts()

Out[2]: Fake    10413
         Real    10387
         Name: label, dtype: int64
```

Model Construction & Discussion

For construction model I have use train dataset and test dataset with the help of Passive Aggressive Classifiers (PAC). For testing the news that it's fake or real I used 33% Test dataset that are contain 20800 article. After I go for vectorization of text and stop al English word like "is", "are", "the" etc.

```
In [3]: x_train,x_test,y_train,y_test=train_test_split(df['text'], df['label'], test_size=0.33, random_state=7, shuffle=True)
tfidf_vectorizer=TfidfVectorizer(stop_words='english', max_df=0.67)

In [4]: vec_train=tfidf_vectorizer.fit_transform(x_train.values.astype('U'))
vec_test=tfidf_vectorizer.transform(x_test.values.astype('U'))
```

Passive Aggressive Algorithm:-

Passive-Aggressive algorithms are generally used for large-scale learning. In machine learning algorithms, the input data comes in sequential order and the machine learning model is updated step-by-step, as opposed to batch learning, where the entire training dataset is used at once. This is very useful in situations where there is a huge amount of data and it is computationally infeasible to train the entire dataset because of the sheer size of the data.

There are some important parameter of PAC that how it work-

- † C : This is the regularization parameter, and denotes the penalization the model will make on an incorrect prediction
- † max_iter : The maximum number of iterations the model makes over the training data.
- † tol : The stopping criterion.

In my model max_iter I have taken max_iter =150. It's default we can take any value whatever we want.

After doing this I predict with this (PAC) model that how much Accuracy It is give For checking accuracy I see that it's give good accuracy(96.24%). It mean that almost 96% I

```
In [5]: pac=PassiveAggressiveClassifier(max_iter=150)
pac.fit(vec_train,y_train)

Out[5]: PassiveAggressiveClassifier(C=1.0, average=False, class_weight=None,
                                     early_stopping=False, fit_intercept=True,
                                     loss='hinge', max_iter=150, n_iter_no_change=5,
                                     n_jobs=None, random_state=None, shuffle=True,
                                     tol=0.001, validation_fraction=0.1, verbose=0,
                                     warm_start=False)
```

predict correctly that given article is fake or real.

Confusion Matrix:-

A confusion matrix is a table that is often used to describe the performance of a classification model on a set of

test data for which the true values are known. In my model result of confusion matrix is

```
array([[3304, 133],  
       [125, 3302]])
```

It mean that 3304 times it predict real and 3302 times it predict fake.

If I talk about K fold Accuracy than it's give 96.20% that is almost near to PAC accuracy. After I take some separate dataset name as "True" and "Fake" that already have separate true and fake articles like

If I show df_true then output is

		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 diplomati...	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
...	
21412		'Fully committed' NATO backs new U.S. approach...	BRUSSELS (Reuters) - NATO allies on Tuesday we...	worldnews	August 22, 2017	Real
21413		LexisNexis withdrew two products from Chinese ...	LONDON (Reuters) - LexisNexis, a provider of l...	worldnews	August 22, 2017	Real
21414		Minsk cultural hub becomes haven from authorities	MINSK (Reuters) - In the shadow of disused Sov...	worldnews	August 22, 2017	Real
21415		Vatican upbeat on possibility of Pope Francis ...	MOSCOW (Reuters) - Vatican Secretary of State ...	worldnews	August 22, 2017	Real
21416		Indonesia to buy \$1.14 billion worth of Russia...	JAKARTA (Reuters) - Indonesia will buy 11 Sukh...	worldnews	August 22, 2017	Real

21417 rows x 5 columns

And df_fake show the output is

	title	text	subject	date	label
0	Donald Trump Sends Out Embarrassing New Year'...	Donald Trump just couldn't wish all Americans ...	News	December 31, 2017	Fake
1	Drunk Bragging Trump Staffer Started Russian ...	House Intelligence Committee Chairman Devin Nu...	News	December 31, 2017	Fake
2	Sheriff David Clarke Becomes An Internet Joke...	On Friday, it was revealed that former Milwauke...	News	December 30, 2017	Fake
3	Trump Is So Obsessed He Even Has Obama's Name...	On Christmas day, Donald Trump announced that ...	News	December 29, 2017	Fake
4	Pope Francis Just Called Out Donald Trump Dur...	Pope Francis used his annual Christmas Day mes...	News	December 25, 2017	Fake
...
23476	McPain: John McCain Furious That Iran Treated ...	21st Century Wire says As 21WIRE reported earl...	Middle-east	January 16, 2016	Fake
23477	JUSTICE? Yahoo Settles E-mail Privacy Class-ac...	21st Century Wire says It's a familiar theme. ...	Middle-east	January 16, 2016	Fake
23478	Sunnistan: US and Allied 'Safe Zone' Plan to T...	Patrick Henningsen 21st Century WireRemember ...	Middle-east	January 15, 2016	Fake
23479	How to Blow \$700 Million: Al Jazeera America F...	21st Century Wire says Al Jazeera America will...	Middle-east	January 14, 2016	Fake
23480	10 U.S. Navy Sailors Held by Iranian Military ...	21st Century Wire says As 21WIRE predicted in ...	Middle-east	January 12, 2016	Fake

23481 rows x 5 columns

After I make final model that predict using our model that news are fake or real. for this I build logic that give 69.83% time give correct and 69.07% time tell that given article are fake.

```
In [13]: sum([1 if findlabel((df_true['text'][i]))=='Real' else 0 for i in range(len(df_true['text']))])/df_true['text'].size
```

Out[13]: 0.6983704533781575

```
In [14]: sum([1 if findlabel((df_fake['text'][i]))=='Fake' else 0 for i in range(len(df_fake['text']))])/df_fake['text'].size
```

Out[14]: 0.6907286742472637

And finally I take another dataset (model) and make function fakenews to predict that given news is real or fake. In this model It will take input from user and give the corresponding output either true or fake.

For prediction I take 2 input first “Obama is the president of USA in 2018” and second is “Donald Trump is the president of USA in 2018” then for first input it give “false” output and for second it give “true”.

```
In [28]: import pickle
s1=input("Enter the news that you want to predict: ")
print("You entered: "+str(s1))

def fakenews(s1):
    load=pickle.load(open('C:/Users/Sonu/fake_news_detector/model.sav','rb'))
    prediction=load.predict([s1])
    prob=load.predict_proba([s1])
    return (print("The given statement is",prediction[0]),
            print("the truth probability score is",prob[0][1]))

if __name__=='__main__':
    fakenews(s1)
```

```
Enter the news that you want to predict: Obama is the president of USA in 2018\
You entered: Obama is the president of USA in 2018\
The given statement is False
the truth probability score is 0.349987998991195
```

```
Enter the news that you want to predict: Donald trump is the president of USA in 2018\
You entered: Donald trump is the president of USA in 2018\
The given statement is True
the truth probability score is 0.5958993316009602
```

Conclusion

The problems related to fake news have increased considerably in the last years, particularly in the political sector. In this model, I presented a fake news detection model using Machine learning and Passive Aggressive Classifiers. I have incorporated different metadata (text, author, and title) to perform the fake news detection. As a result, the model has reached its highest accuracy with the text and author input . The highest PAC accuracy score is 96.24%. In my future work, I will run our model on large datasets and incorporate multiple metadata in different context to create a complete picture of fake news detection performance.

References

- 1) Kaggle.<https://www.kaggle.com/c/fake-news>
- 2) Kwon, S., Cha, M., Jung, K., Chen, W., Wang, Y. (2013, December). Prominent features of rumor propagation in online social media. In 2013 IEEE 13th International Conference on Data Mining (pp. 1103-1108). IEEE
- 3) L. Luceri, A. Vancheri, T. Braun, and S. Giordano, On the social influence in human behavior: Physical, homophily, and social communities, in Proceedings of the Sixth International Conference on Complex Networks and Their Applications, 2017., 2017, pp. 856–868
- 4) Geeksforgeeks.<https://www.geeksforgeeks.com/passive-aggressive-classifiers/>
- 5) AHMED, Hadeer, TRAORE, Issa, et SAAD, Sherif. Detection of online fake news using n-gram analysis and machine learning techniques. In : International Conference on Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments. Springer, Cham, 2017. p. 127-138