

Interpretability of Fake News Detection Model

Likhitha Pulluru¹, Laxmi Shravani Mamidala¹,

Dr. Ramanathan Subramanian Prof², and Abhinav Dhall Prof²

¹ *INSOFE*

² *IIT Ropar*

ABSTRACT

The authenticity of information has been a longstanding issue with its potential to impact millions of users in the blink of an eye. Recent years has seen a growth in development of fake news detection models. Model Interpretability especially in NLP domain is still challenging, yet it helps in adopting models to various domains. In this study we tried post hoc interpretation techniques like local and global interpretations using LIME and SHAP, Topic modeling and Keyword extraction techniques. We identified these are simple yet powerful techniques to better understand the tagging behavior of fake news detection models used in this paper. With the help of these interpretation and data augmentation techniques we measured model robustness and identified that models built on ML algorithms are not robust to covariance shift in input data. Also, we tried to derive some of the characteristics that better represents style of fake tweets.

1. INTRODUCTION

Any untrue information disguised as a credible news source is termed as fake news. Irony is fake news doesn't carry any pre-defined characteristics which makes it highly difficult to classify. Adopting models across different domains is not easy as it shows high variability across domains. Recent times fake news related to COVID-19 is massive. We have made our study on COVID related tweets data taken from Fighting an Infodemic: COVID-19 Fake News Data set paper [1] In this article, we explored dataset taken from above paper and built machine learning models using SVM, Random Forest and Gradient Boost. All these algorithms gave F1-score in range of 91% – 93%. Then we moved on to post hoc model interpretations using LIME and SHAP which helped in identifying contribution of each word in classifying the tweets. We performed topic modeling to understand the difference in origin of domains for fake and real tweets. Finally, we used keywords extraction technique using pre-built python libraries and observed high scores are result of raw keywords presence in different classes. This observation gave us a stand of, using simple ML algorithms for fake news detection problems are not reliable even with high F1-scores. To better support this statement, we measured model robustness using data augmentation techniques and by inducing co-variance shift in input data. Finally, we put forward scope of future work for making fake news detection models robust and adoptable to different domains. The code is available at [https://github.com/](https://github.com/likhitha79/Fake_News_Detection_Model_Interpretability)

[likhitha79/Fake_News_Detection_Model_Interpretability](https://github.com/likhitha79/Fake_News_Detection_Model_Interpretability)

1.1 Literature Survey

Human Fact checkers are considered to be gold standard most of them time. However, studies in social psychology and communications have demonstrated that human ability to detect deception is only slightly better than chance. Typical accuracy rates are in the range of 55%-58%. Manual Expert based fact checking websites (like PolitiFact, snopes, TextThresher etc) have emerged to classify fake news. Nevertheless, manual fact checking does not scale well with the volume of data getting generated on social media.

Can Data Science Identify fake news?

Full Fact is a UK based fact checking outfit recently received grants from Google for their work in this area. Interesting point is they did not choose traditional path of NLP which is similar to spam detection. They built a Watson like platform that can parse facts floating around the world as unstructured data and using it as base to classify news. This is like loading the system with huge volume of curated known facts and then comparing new material using the logic of QAM's (Question Answering Machines). In the advent of building efficient fake news detection systems models incorporated are becoming more complicated. Interpreting these models is need of the hour. Recently attention mechanism has been widely used in var-

ious NLP tasks and their attention scores have been used as a technique for interpretation. However this interpretation mechanism is still controversial [2]

Using cloning technique for model interpretation where in some dimensions of the embedding vector will be erased and the result is analyzed. However, erased vectors will be replaced with zero. Such erasure scheme will push out some data from training data distribution thereby resulting in inaccurate interpretation. In this paper **Interpretation of NLP models through input marginalization** [3] they proposed a new technique by marginalizing each token out to mitigate OOD problem of existing erasure problem.

All the above studies are about NLP models in general and we didn't find much content pertaining to interpretation of fake news detection models. This gave us motivation to explore much deeper into interpretability part.

2. DATASET EXPLORATION

In Fighting an Infodemic paper **Fighting an Infodemic: COVID-19 Fake News Dataset** [1] authors have curated dataset by collecting tweets related to COVID from various social media platforms. Authenticity is checked using multiple fact checking websites like Politifact2, NewsChecker3, Boomlive4

2.1 Basic Description of data set

This dataset has train, validation, and test sets.

Fig. 2..1. Data set Information

train_data.info()					test_data.info()				
<class 'pandas.core.frame.DataFrame'> RangeIndex: 6420 entries, 0 to 6419 Data columns (total 3 columns): # Column Non-Null Count Dtype --- --- --- --- --- 0 id 6420 non-null int64 1 tweet 6420 non-null object 2 label 6420 non-null object dtypes: int64(1), object(2)					<class 'pandas.core.frame.DataFrame'> RangeIndex: 2140 entries, 0 to 2139 Data columns (total 2 columns): # Column Non-Null Count Dtype --- --- --- --- --- 0 id 2140 non-null int64 1 tweet 2140 non-null object dtypes: int64(1), object(1)				
val_data.info()									
<class 'pandas.core.frame.DataFrame'> RangeIndex: 2140 entries, 0 to 2139 Data columns (total 3 columns): # Column Non-Null Count Dtype --- --- --- --- --- 0 id 2140 non-null int64 1 tweet 2140 non-null object 2 label 2140 non-null object dtypes: int64(1), object(2)									

Data Dictionary:

Id: Unique identifier for each message

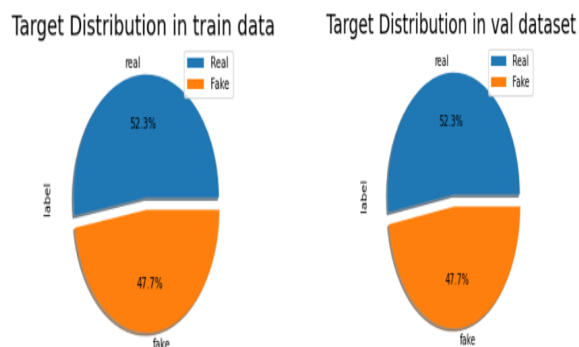
Tweet: Actual Message

Label: Tag indicating fake or real news

2.1.1 Distribution of target variable

Let's look at distribution of fake/ real tweets in train and validation datasets.

Fig. 2..2. Distribution of Target variable



Data is looking quite balanced in both train and validation datasets.

2.2 Characteristics of fake tweets

Polarity and subjectivity of the data :

Text polarity describes whether it is a positive(+1), negative(-1) or neutral(0) statement varying from -1 to +1. Whereas text subjectivity is a measure of how subjective or objective the statement is, for e.g. text with more of opinions, emotions, judgements carry high subjectivity score rather text with factual information. We have used **TextBlob** for calculating polarity and subjectivity of tweets

Fig. 2..3. Polarity Distribution in Real Tweets

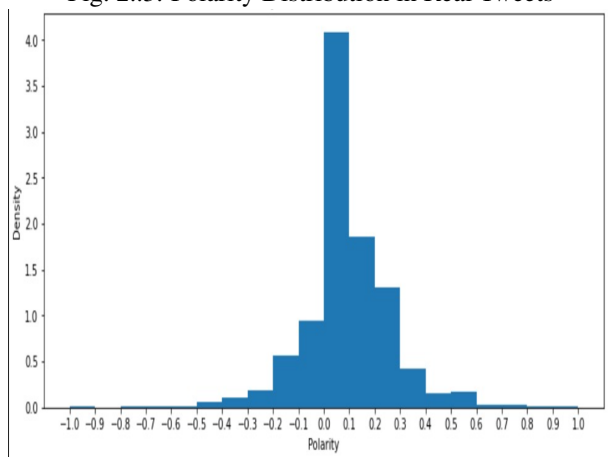
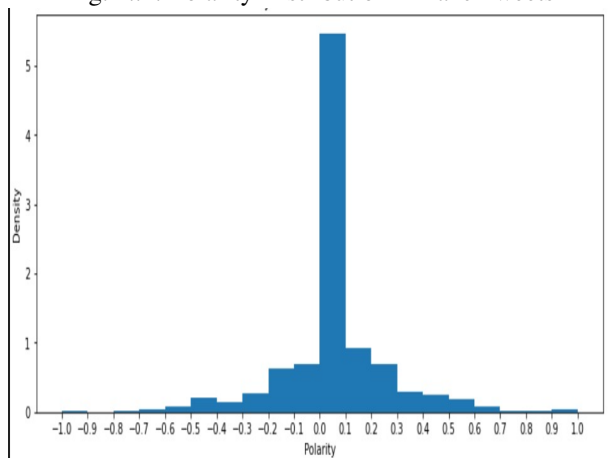


Fig. 2..4. Polarity Distribution in Fake Tweets



- Fake tweets are little negative skewed compared to real tweets

Polarity:

Fig. 2..5. Z-Score Distribution

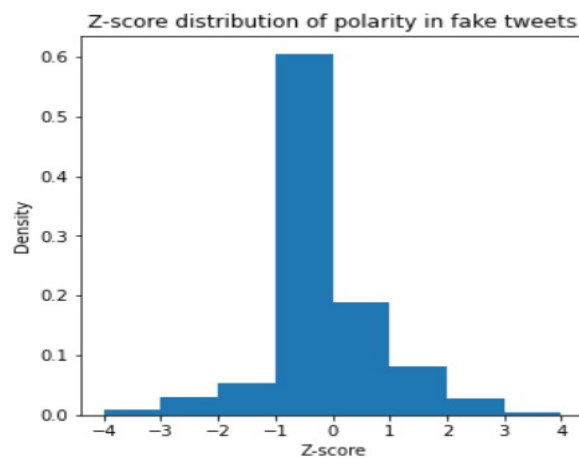
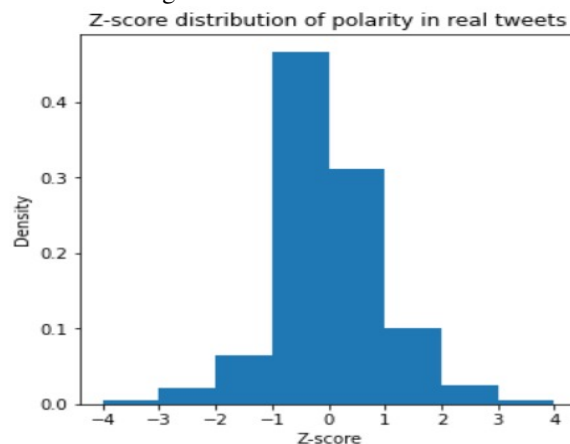


Fig. 2..6. Real tweets polarity

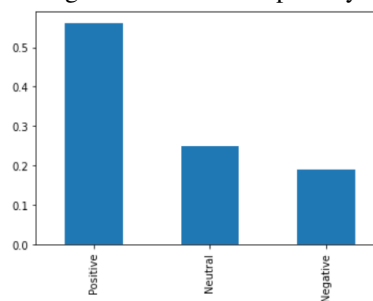
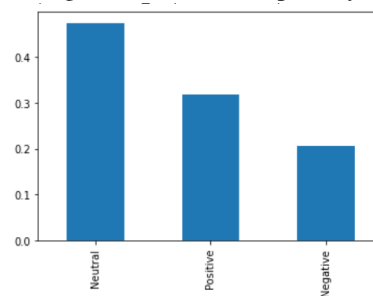


Fig. 2..7. Fake tweets polarity



- Majority of fake tweets are carrying neutral sign whereas real tweets are mostly in positive tone

Table. 1. Mean and variance of polarity values in real and fake tweets

	Mean (Polarity)	Variance (Polarity)
Real Tweets	0.074	0.031
Fake Tweets	0.025	0.042

Table. 2. Mean and variance of Subjectivity values in real and fake tweets

- It seems subjectivity of fake tweets are centered around 0 indicating most of them are factual information rather opinions

	Mean (Subjectivity)	Variance (Subjectivity)
Real Tweets	0.410	0.07
Fake Tweets	0.263	0.08

Subjectivity:

Fig. 2..8. Subjectivity Distribution in Real Tweets

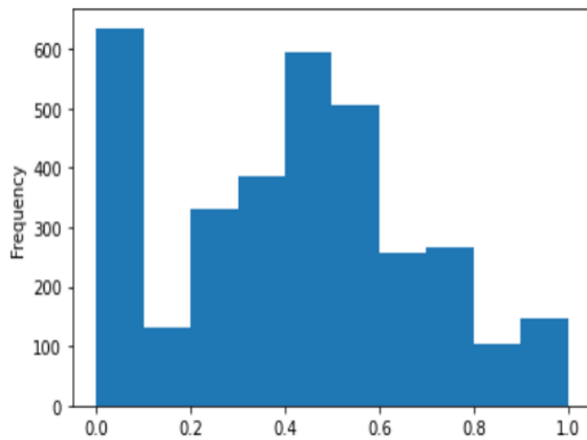
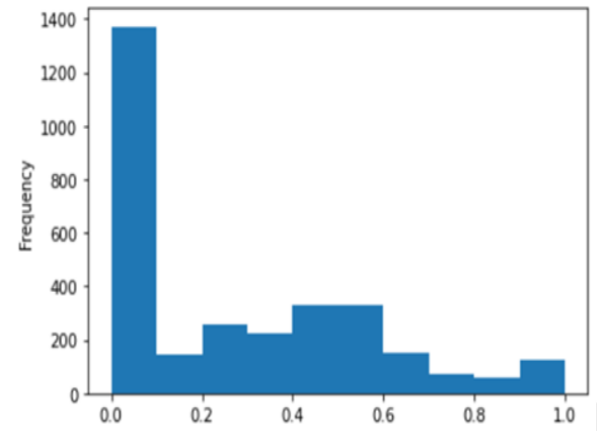


Fig. 2..9. Subjectivity Distribution in Fake Tweets



Distribution of special characters in real and fake tweets:

Real Tweet Data Distrubution

Max_Spl

```

/      1979
@      358
#      329
:      178
%       36
&       35
(       33
?        5
*        5
_        3
!        3
~        2
)        2

```

Name: id, dtype: int64

Fake Tweet Data Distrubution

Max_Spl

```

/      820
#      216
?      207
@      105
(       98
_       72
:       52
!       50
%       24
$       16
*       13
&       10
~        1

```

Name: id, dtype: int64

We can observe that fake tweets carry more “?” characters than real tweets.

Fig. 2..10. Special Characters Distribution in Tweets

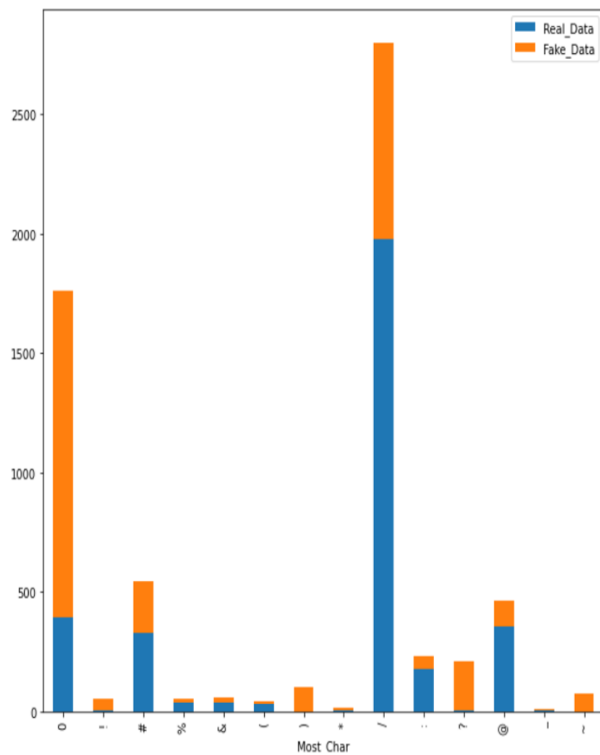
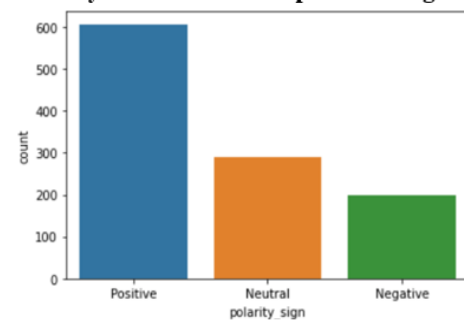


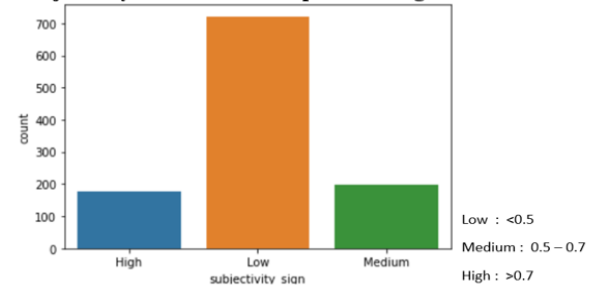
Fig. 2..11. Wordcloud of hashtags appeared only in Real tweets :



Polarity distribution of top 15 hashtags in real tweets

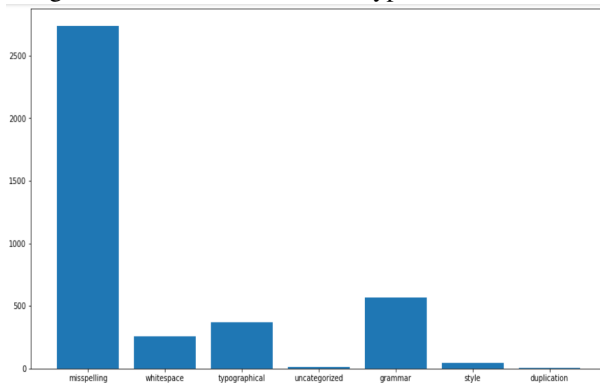


Subjectivity distribution of top 15 hashtags in real tweets



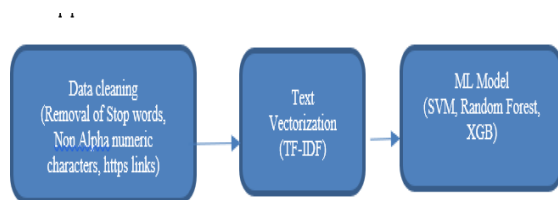
Hashtag information in tweets may have some information to classify between real/fake tweets. Let's look at their distribution.

Fig. 2..16. Distribution of error types in fake tweets :



3. MODEL BUILDING

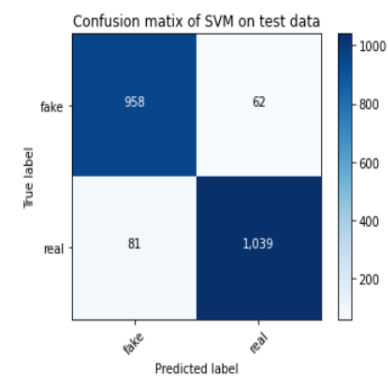
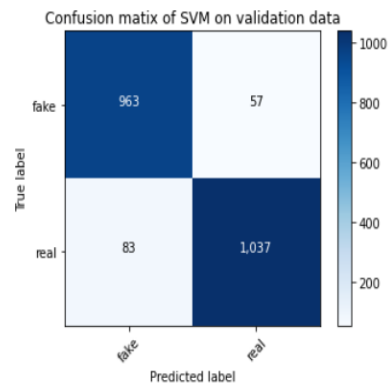
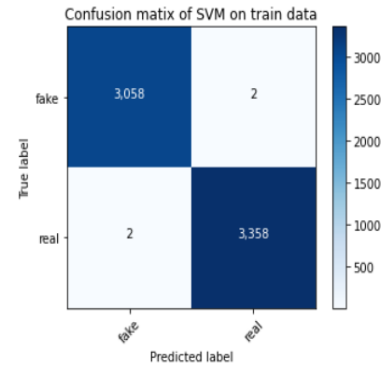
Here is the pipeline used for fake news detection model



As mentioned in Fighting an Infodemic paper we got better results with SVM. Below are the metrics for validation and test sets.

3.1 Evaluation metrics of SVM model

We have explored with SVM, random forest and gradient boosting techniques. However, SVM got better results.



Evaluation metrics on validation dataset:

Model	Accuracy	Precision	Recall	F1
SVM	93.78	93.78	93.7	93.74
Random Forest	92.3	92.3	92.3	92.3
Decision Trees	88.03	88.03	88.03	88.03

Models are giving great scores on both validation and test

sets. Let's look at model interpretability to better understand model's behavior.

4. Model Interpretability

4.1 Using LIME/SHAP

LIME LIME method focus on interpreting locally instead of providing global interpretation. It zooms in a data-point in a model and then find out which features impacted the model to reach certain conclusion. It doesn't care about ML model being used as it treats every model as a black box.

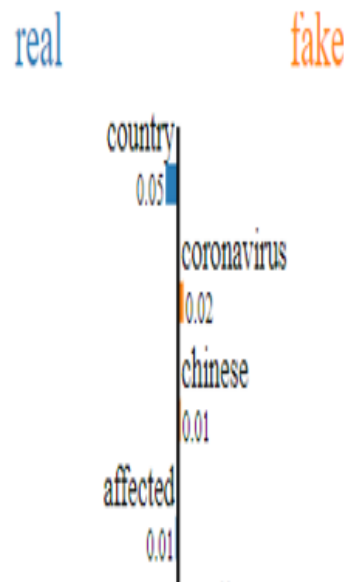
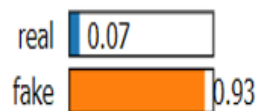
Instance explanations using LIME : (on random validation data)

1. Document id: 0 ; Predicted class = fake ; class: fake

('country', -0.053), ('coronavirus', 0.023), ('Chinese', 0.011), ('affected', -0.010), ('Muslim', 0.001), ('Islam', 0.001), ('covid19', 0.0005), ('realising', 0.0005), ('converting', -0.0004)

Fig. 4..1. Instance Explanation for Document ID :0

Prediction probabilities



Text with highlighted words

chinese converting islam realising muslim affected coronavirus covid19 country

2.Document id: 3; Predicted class = fake; True class: fake

('state', -0.07), ('19', 0.06), ('trump', 0.037), ('president', 0.012), ('donald', 0.011), ('partnership', -0.009), ('pence', -0.007), ('praises', -0.005), ('seamless', -0.0047), ('speech', -0.0038), ('rnc2020', -0.0036), ('covid', -0.0034), ('rnc', -0.002834487786151808), ('mike', -0.0024), ('governors', 0.0019)

Fig. 4..2. Instance Explanation for Document ID :3

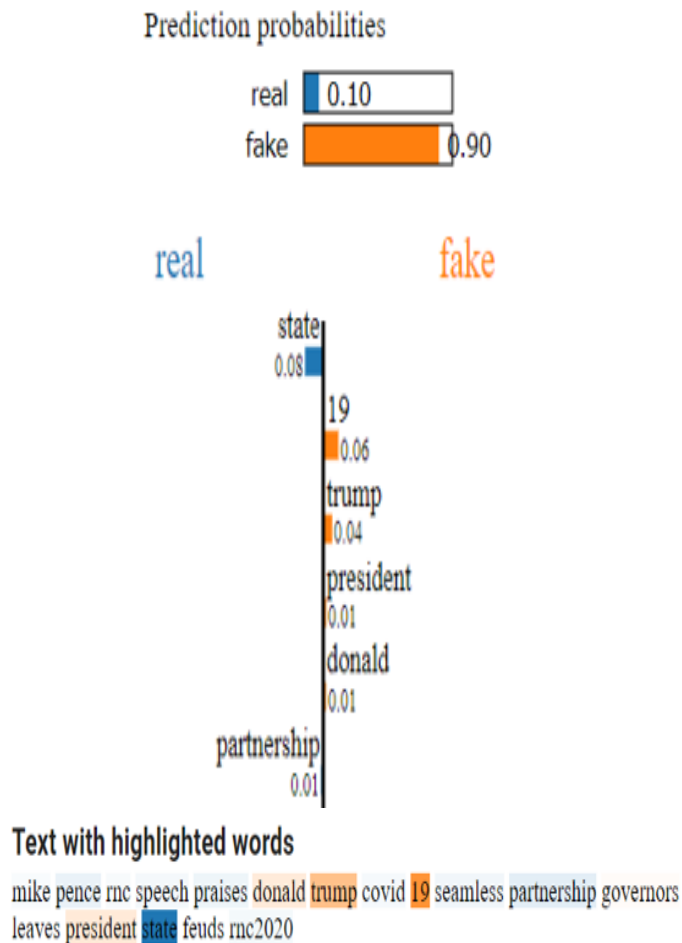
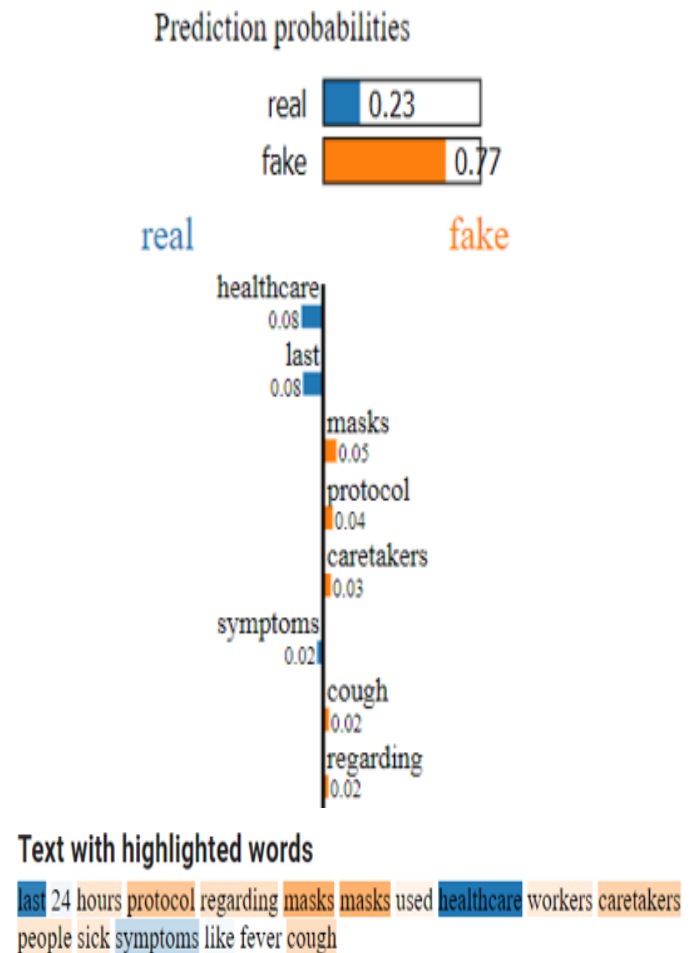


Fig. 4..3. Instance Explanation for Document ID :1010



3.Document id: 1010; Predicted class = fake;True class: fake

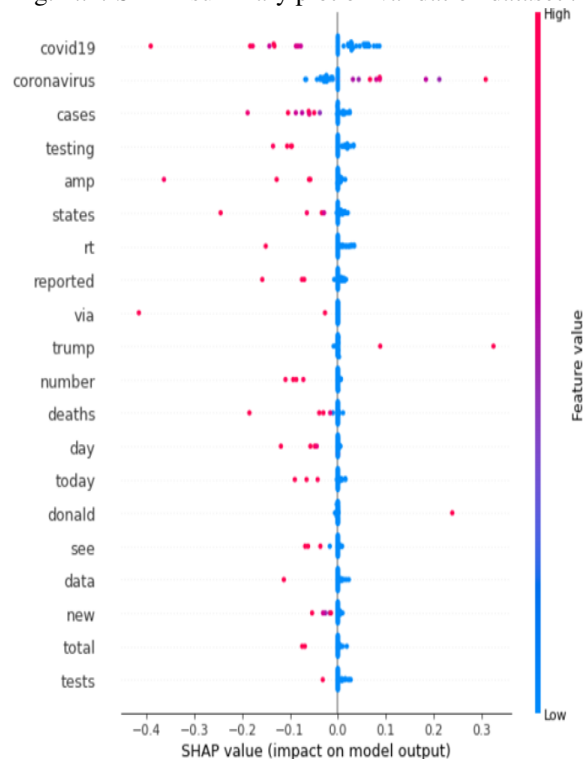
('healthcare',-0.084),('last',-0.078),('masks',0.05),('protocol',0.03),('caretakers',0.02),('symptoms',-0.022),('cough',0.02)

As detailed above, Lime gives out only instance explanations. We have explored with SHAP to get global interpretations.

SHAP

SHAP(Shapley Additive Explanations), It's an average of marginal contributions across all permutations.It can provide both local and global interpretations

Fig. 4.4. SHAP summary plot on validation dataset :



4.2 Topic Modeling

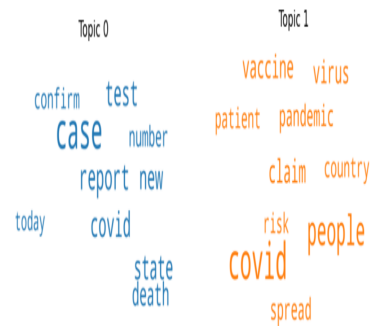
It's an unsupervised NLP technique used to represent a text document with the help of certain topics, that can best explain the underlying information in documents. LDA (Latent Dirichlet Allocation) is one of the powerful algorithms used for topic modelling.

Basic working of LDA :

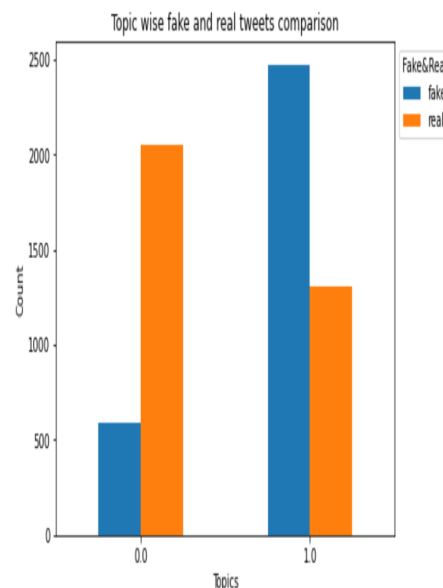


Each document is modeled as multinomial distribution of topics and each topic is modeled as a multinomial distribution of words. It works based on the assumption that every document explains certain topics and those topics generate words based on their probability distribution.

Topic modeling on our fake news dataset :



As dataset is curated specially on COVID domain, we couldn't get clear demarcation between topics derived. we can observe topic 0 is more about covid cases and reports however, topic 1 is more of vaccination and spread. Let's look at the distribution of fake and real tweets across the topics formed.



4.3 Key Words Extraction

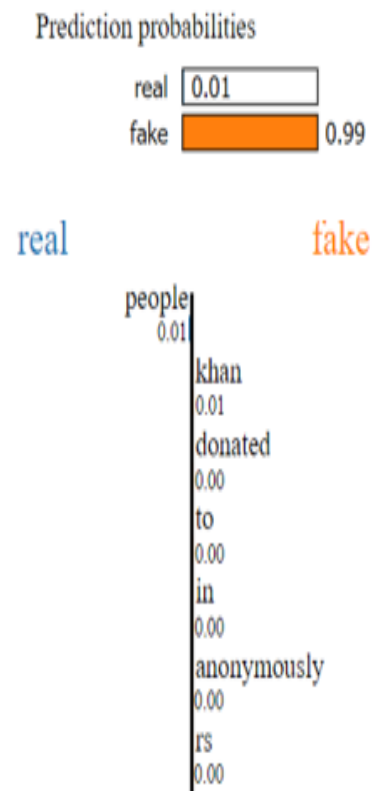
It's a text analysis technique that automatically extracts the most used and important words from a text. It uses word collocations (N-grams) and TF-IDF scores to determine importance of word. Below are the pre built python libraries most used for keyword extraction: 1.Spacy 2.Yake 3.Rake – NLTK 4.Genism **Sample fake tweets :**

• Bollywood actor **Aamir Khan** has anonymously donated Rs 15000 to people living in a slum. === **Fake**

- Bollywood actor **Aamir Khan** actually found a unique way to help the poor people. He filled a truck carrying a bag of flour and went to a locality and called people from their homes, and it was said that only 1 kg. You will get flour. **Fake**
- Film star **Aamir Khan** distributed 15 thousand rupees to the poor in flour bags. **Fake**
- **Aamir Khan** Donate 250 Cr. In PM Relief Cares Fund **Fake**
- Video of an elderly woman struggling to breathe lying over what seems to be a plastic bag. According to the post she was inside a bag at a hospital morgue and was rescued by her relatives. Politically Correct Woman (Almost) Uses Pandemic as Excuse Not to Reuse **Plastic Bag**. **Fake**
- Video shows muslim women spitting in **plastic bags** and throwing them into the houses to spread coronavirus. **Fake**

Document id: 406; Predicted class = fake; True class: fake

Fig. 4..5. Instance Explanation for Document ID :406



Text with highlighted words

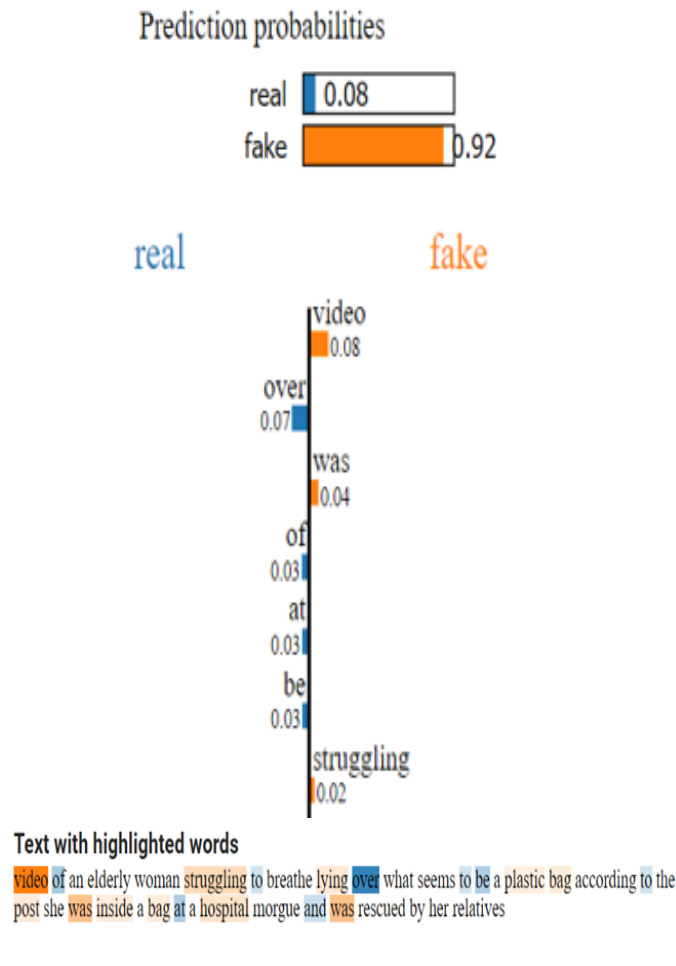
bollywoof actor aamir **khan** has anonymously donated rs 15000 to **people** living in a slum

Document id: 503; Predicted class = fake; True class: fake

('people', -0.012), ('khan', 0.006), ('donated', 0.004), ('to', 0.004),
 ('in', 0.002), ('anonymously', 0.001), ('rs', 0.0012),
 ('aamir', 0.001), ('actor', 0.0009), ('bollywood', 0.0009),
 ('slum', 0.0009), ('15000', 0.0007),
 ('has', -0.0006), ('living', 0.0006)

('video', 0.07), ('over', -0.070), ('was', 0.037), ('of', -0.02),
 ('at', -0.02), ('be', -0.02), ('struggling', 0.02), ('and', -0.018),
 ('hospital', 0.01), ('to', -0.016), ('lying', 0.013), ('post', 0.01), ('bag', 0.01),
 ('inside', 0.011), ('plastic', 0.01)

Fig. 4..6. Instance Explanation for Document ID :503



Also we have tested predicting tweets which contains these impactful words, irrespective of the credibility of the news if tweet contains these words like Aamir Khan, Plastic Bags it is giving fake tag for all the tweets.

Sample Real tweets :

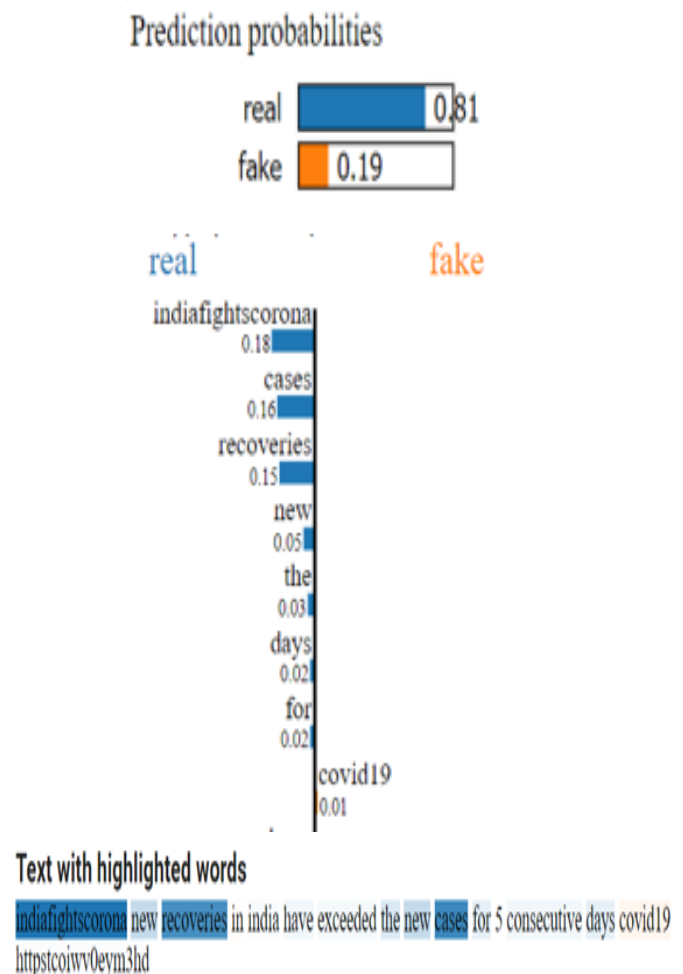
- **IndiaFightsCorona:** 1098621 tests were conducted in the last 24 hours testifying the enlarged testing capacity in the country.=== **Real**
- **IndiaFightsCorona:** Maharashtra Karnataka and Andhra Pradesh. Together with the States of Uttar Pradesh and Tamil Nadu these 5 states contribute nearly 60% of the total active cases.=== **Real**
- **IndiaFightsCorona:** Guidelines for phased re-opening Unlock4 StaySafe IndiaWillWin.=== **Real**
- **CoronaVirus Updates IndiaFightsCorona:** Total COVID19 recovered cases has touched another high of 1480884 today. This is 2.3 times the number of active cases (628747 today).

Case Fatality Rate has further slumped to 2.01% === **Real**

Document id: 45 ;Predicted class = real; True class: real

('indiafightscorona',-0.183),('cases',-0.160),('recoveries',-0.150),('new',-0.04),('the',-0.029),('days',-0.02),('for',-0.019),('covid19',0.012),('have',-0.011),('consecutive',-0.011),('exceeded',-0.008),('india',-0.005),('in',0.004),('httpstcoiww0eym3hd',-0.003)

Fig. 4..7. Instance Explanation for Document ID :45



Also we have made a test run using all tweets contains "IndiaFightsCorona". Again without considering much about credibility model gave a Real tag for almost all the tweets. It failed to tag some of the fake tweets containing this keyword.

Here is the interesting observation, certain keywords are

highly contributing for label tagging as real or fake. To explore more about this we selected few tweets and extracted keywords out of it using **Yake** library as it provides option for n-grams selection and duplication threshold. Then for each keyword in tweet we have calculated number of times that keyword occurred in whole corpus of real and fake tweets and formed a dictionary with real and fake label counter. For instance, In below fake tweets, if we see total label counter for all keywords they have occurred 28 times in real tweets and 56 times in fake tweets. Thus, making these tweets as fake one's.

Aamir Khan

```
[137] dict1_counter
{
  'aamir': [0, 5],
  'aamir khan': [0, 5],
  'actor': [15, 14],
  'actor aamir': [0, 3],
  'anonymously donated': [0, 1],
  'bollywood actor': [0, 1],
  'donated': [0, 5],
  'khan': [7, 17],
  'people living': [3, 4],
  'slum': [3, 1]}

[140] total_label_counter(dict1_counter)
(28, 56)
```

Plastic Bag

```
[140] dict2_counter
{
  'bag': [0, 16],
  'breathe lying': [0, 1],
  'elderly woman': [0, 1],
  'hospital morgue': [0, 1],
  'lying': [18, 26],
  'plastic bag': [0, 3],
  'relatives': [1, 3],
  'video': [9, 189],
  'woman': [17, 34],
  'woman struggling': [0, 1]}

[141] total_label_counter(dict2_counter)
(45, 275)
```

Similarly, for some real tweets we got real value counter more than fake one's as below

IndiaFightsCorona

```
[159] dict4_counter = label_counter(train_data,myDict4)
dict4_counter
{
  'cases': [976, 131],
  'consecutive days': [3, 0],
  'days': [224, 70],
  'exceeded': [12, 0],
  'india': [563, 255],
  'indiafightscorona': [319, 3],
  'recoveries': [99, 4]}

[160] total_label_counter(dict4_counter)
(2196, 463)
```

Handwashing

```
[163] dict5_counter = label_counter(train_data,myDict5)
dict5_counter
{
  'dyk': [9, 2],
  'dyk handwashing': [1, 0],
  'handwashing remains': [1, 0],
  'prevent': [83, 70],
  'remains': [89, 2],
  'spread': [192, 79],
  'surgeongeneral dyk': [2, 0],
  'things': [10, 11],
  'viruses': [5, 26]}

[164] total_label_counter(dict5_counter)
(392, 190)
```

We can understand clearly from above examples that model

predictions are heavily driven by keywords and their appearance in fake and real tweets.

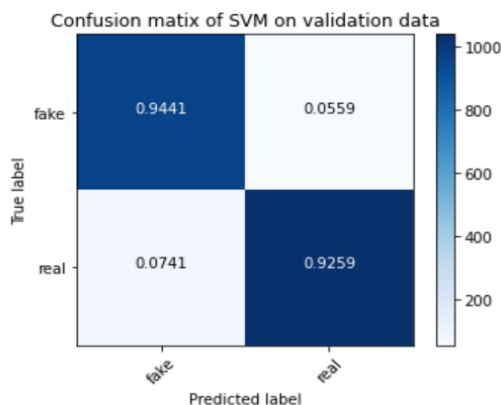
5. Measuring Model Robustness

These are the state of art data augmentation techniques generally used to augment data. However, we have used them to check if model interpretability depends on keywords extraction and their appearance in fake/real tweets. • Random synonym replacement • Random Insertion • Random swap

5.1 Random Swap

Here we have swapped all the words (formed jumbled sentence) in a tweet. As the model we have built in section 3 doesn't bother about context and order of words its performance will not get impact.

```
Accuracy : 0.9345794392523364
Precision : 0.9348008619335585
Recall : 0.9345794392523364
F1 : 0.9345519215989282
```



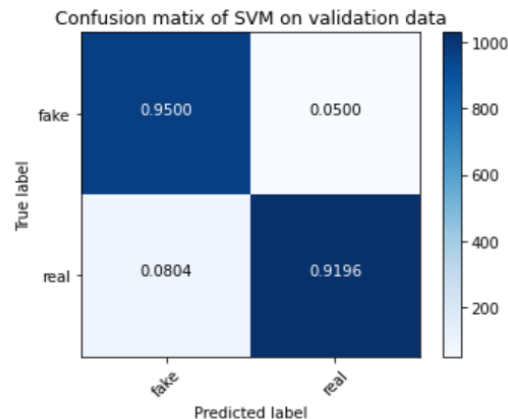
As expected we got F1 score of 93.45% (similar to base model) for swapped tweets.

5.2 Random Insertion

In here we have randomly inserted synonyms in a tweet. As model built highly depends on raw keywords (as observed in section 4.3) this insertion also doesn't give much impact on the model performance. Here are the evaluation metric details on validation dataset after random insertion of syn-

onyms.

```
Accuracy : 0.9341121495327103
Precision : 0.9346653871829106
Recall : 0.9341121495327103
F1 : 0.9340778944319391
```

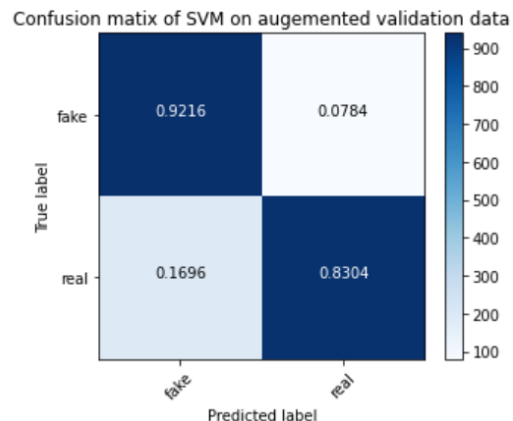


We got F1-score of 92.7 thus showing very less impact on model performance.

5.3 Random synonym replacement

Here we have replaced randomly selected keywords (taken from keywords extraction using Yake) with their synonyms taken from word net library. Also replaced nouns with static text as there will be no synonyms for nouns. Below are the evaluation metric details. **Sample tweet replaced with synonyms:**

```
Accuracy : 0.8738317757009346
Precision : 0.8785202164978141
Recall : 0.8738317757009346
F1 : 0.873862081421369
```

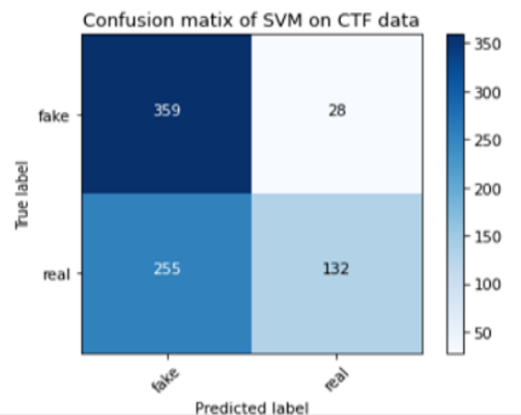


As model highly depends on raw keywords as shown in section 4.3, synonym replacement will have high impact on the model's performance thus bringing down F1 score to 75%..

5.4 Co-variance Shift

This refers to the change in the distribution of input data. This is very common in real world as data can be collected from different sources. Here we have tested model robustness to shift in input by taking data set from same domain (COVID-19) but from different sources. As our base model uses raw keywords for tagging the tweet, co-variance shift highly impacts the model performance and it got down F1 score to 66.8%.

```
Accuracy : 0.6343669250645995
Precision : 0.80639518191348
Recall : 0.6343669250645995
F1 : 0.6687762209211128
```



By looking at the model robustness measures, it's clear that fake new detection models built on ML algorithms cannot be reliable even after having high F1-scores.

6. CONCLUSION AND FUTURE SCOPE

In this article, we explored data-set given and tried identifying characteristics of fake news. We have looked at various model interpretation techniques to understand predictions of fake news detection models built using ML algorithms. Also, we have seen performance measuring techniques to check model robustness. Detecting fake news is a challenging task as it's characteristics can vary based on various factors like originating source and domain. However, considering style of writing may give some sort of clue in better

differentiating real/fake tweets. Also, considering state of art algorithms in NLP for model building can help in giving out better predictions based on actual context.

BIBLIOGRAPHY

- [1] P. Patwa, S. Sharma, S. Pykl, V. Guptha, G. Kumari, M. S. Akhtar, A. Ekbal, A. Das, and T. Chakraborty, "Fighting an infodemic: Covid-19 fake news dataset," in *International Workshop on Combating On line Ho st ile Posts in Regional Languages dur ing Emerge ncy Si tuation*, pp. 21–29, Springer, 2021.
- [2] S. Jain and B. C. Wallace, "Attention is not explanation," *arXiv preprint arXiv:1902.10186*, 2019.
- [3] S. Kim, J. Yi, E. Kim, and S. Yoon, "Interpretation of nlp models through input marginalization," *arXiv preprint arXiv:2010.13984*, 2020.