A DEEP MULTIMODAL MODEL FOR FAKE TWEET DETECTION USING CONTENT BASED, METADATA AND DERIVED FEATURES

VAISHALI VAIBHAV HIRLEKAR

Sir Padmapat Singhania University Udaipur, Rajasthan, India. Email: vaishali.hirlekar@spsu.ac.in

ARUN KUMAR

Sir Padmapat Singhania University Udaipur, Rajasthan, India. Email: arun.kumar@spsu.ac.in

Abstract

Social media has grown in importance over the last few decades because it allows people from all over the world to stay connected, however, it has become a breeding ground for misinformation during different public event as well as in the case of COVID-19 pandemic. Detection of fake information on social media has been technically challenging as it necessitates time-consuming evidence gathering and meticulous fact checking. There are generally three widely accepted characteristics of fake tweet: tweet content and associated features/metadata and credibility of the source. Using these characteristics, we propose a novel Fake Tweet Detection model in this paper that can identify whether given tweet is real or fake. The primary objective of the research is to identify fake tweets, and experimentation is carried out precisely utilizing COVID-19 fake tweets as a case study. We have used an ensemble model composed of tweet text features, metadata features and derived features present in the tweets. In order to assess the efficacy of the proposed methodology, we assessed our model using the COVID-19 dataset and achieved 99.42 F1-Score

Keywords: Social media, Fake Tweet, Deep learning, Covid 19

1. INTRODUCTION

Social media is becoming more relevant and widely used than ever thanks to technological advancements. However, social media platforms impacts the quick generation of a substantial amount of information to distribute fake and unreliable information. It can be challenging to identify fake news based alone on news content because it is intentionally created to mislead a wide variety of readers.

Fake news is defined in a number of ways. An article of news that is purposefully and demonstrably untrue is known as fake news ((Shu et al.,2017; Zhang & Ghorbani, 2019; Vishwakarma & Jain, 2020). The phrase "fake news" is used to describe misleading information that appears in mainstream media. (Bondielli & Marcelloni, 2019)

Fake news is a term that is used to describe a variety of ideas, including rumour and misinformation. (Lazer et al.,2018 and Varshney & Vishwakarma, 2021a). According to another definition (Meel & Vishwakarma,2020)., fake

news is a type of misleading information released under the guise of being legitimate news commonly spread through news outlets or the internet with a goal to gain politically or financially, boost reading, and prejudice public opinion. In (Rubin et al., 2015), the authors made distinctions between different types of fake news, such as severe fabrications, massive hoaxes, and hilarious fakes. This fake news has become a breeding ground for misinformation during different public event, political events as well as in the case of COVID-19 pandemic.

Currently, the COVID-19 epidemic is spreading over the globe at an alarming rate. Many believe it to be the greatest global health catastrophe of our time. We're not just fighting an epidemic; we're fighting an infodemic, said Tedros Ghebrevesus, Director-General of the World Health Organization (WHO) in mid-February in Munich, Germany, to a group of security and foreign policy professionals, referring to fake news that ---spreads faster and more easily than this virus. The United Nations Department of Global Communications (DGC), UN.org, 28 Mar.2020. Misinformation about COVID-19 seems to be circulating swiftly on social media. Other epidemics, including the current Ebola and Zika outbreaks, have shown similar characteristics. (Shahi et al., 2021) According to UNESCO, "during this coronavirus pandemic, fake news is putting lives in danger." Fake news, ranging from theories about the origin of the virus to bogus prevention and cures, is rapidly spreading without any valid evidence.(Cinelli et al., 2020; Kouzy et al., 2020; Gallotti et al., 2020; Singh et al., 2020 and Yang et al., 2020) These studies focused at the magnitude and propogation of misinformation concerning the COVID-19 pandemic on Twitter. These studies used a manually annotated subset of Twitter data (Kouzy et al., 2020), with the major issue being that many of them focused on a modest subset of claims (Singh et al., 2020). Other studies (Cinelli et al., 2020; Gallotti et al., 2020 and Yang et al., 2020) used the dependability of the references to automatically detect false information. From the given studies it has been observed that popular source-based methodologies provide extensive analyses of Twitter data (Allen et al., 2020; Zollo et al., 2017; Bovet et al., 2019; Grinberg et al., 2019 and Shao et al., 2018).

In our study, we have tried to identify the Twitter accounts responsible for spreading fake COVID-19 claims and their contents. With these observations, we have employed an approach based on tweet text features, metadata features and derived features present in the tweets. We first give a summary of social media analytics methods that we have used for analysing the tweets related to COVID-19 infodemic. This is the start of a more organised, goaloriented strategy for managing the crisis as it develops and for discovering ways to lessen the damaging impacts of disinformation in crises that arise in the future. Furthermore, we have attempted linguistic features and network features that had an effect on spreading misinformation.

We carried out an exploratory study in order to gain data insights. To detect patterns, spot anomalies, test hypotheses, and verify assumptions, an exploratory data analysis has been performed on the data using summary statistics and graphical representations.

In Section 2, we provide the academic context of our work in the field of misinformation detection and propagation. In Section 3, we elaborate on our data collection and analysis process. We then present the methodology and Implementation details in Section 4 and 5, followed by discussing experimental results in Section 6. Finally, we draw a conclusion in Section 7

2. LITERATURE SURVEY

In this section, we discuss the misinformation and disinformation, the spread of misinformation and the various approaches used for fake tweet detection.

The spread of fake information is a broad definition of misinformation. When fake information is accidentally spread, the term "misinformation" is more frequently employed; by contrast, the term "disinformation" is used to describe fake information that is deliberately spread. In this study, we make no assertions about the intentions of information providers, whether unintentional or malevolent. As a result, regardless of intent, we group fake information rationally.

Studies show that due to the abundance of fake information, people are unable to differ between true and fake news. While numerous fact-checking websites have emerged to help the public in this respect, like Politi Fact and Snopes (Zhang & Ghorbani, 2019), they mainly depend on outside resources to verify the authenticity of news, which poses concerns with efficiency and scalability. Researchers have therefore focused on techniques for automatically identifying fake information. Automatic fake news identification is frequently based on the news's content (Oshikawa et al., 2018). Conversely, social context-based techniques have grown in acceptance and produced remarkable outcomes (Zhou & Zafarani, 2019).

2.1 Detection of fake news based on content

Textual and visual components are typically used to derive content-based features. Expressive writing techniques can be seen in textual features (Ghosh & Shah, 2019; Choudhary, & Arora, 2021; Potthast et al., 2018) besides sentiments and emotions (Zhang & Ghorbani, 2019; Dungs et al., 2018 and Qian et al., 2018). Tensor factorization is mostly used to model and express textual representations. (Hosseinimotlagh & Papalexakis, 2018;

Gupta et al., 2018; Kaliyar et al., 2021) and deep neural systems (Girgis & Gadallah, 2018; Song et al., 2021; Long et al., 2017) which effectively detect fake news. Various aspects of fake news stories are captured using visual elements that are retrieved from photographs and videos (Liang et al., 2015 and Chen et al., 2019).

Author suggested a semi-supervised method for classifying fake news in texts using temporal ensembling convolutional neural networks (Meel & Vishwakarma, 2021a) in this paper. The authors trained the suggested technique by concatenating the feature vectors retrieved after applying various size convolutional filters to the headline and body of news items. The proposed method was successful in correctly differentiating fake news pieces

from authentic material, according to the experimental results analysis.

(Wang et al., 2018) proposed the event adversarial neural network. This model is made up of three major parts: a feature extractor, a fake news detector, and an event discriminator. The textual and visual latent aspects of news are represented as two separate vectors by the two layers of CNNs that make up the multi-model feature extractor. The final multi-model feature representation is created by concatenating the latent feature representations of the text and the image after they have been learned. The fake news detector then determines if the news is authentic or fraudulent using the multimodel feature representation. The event discriminator predicts the label of the event connected to each news article using the multi-model feature representation as input.

Other studies investigated content-based models using a variety of techniques, including attention-residual networks (Chen et al., 2019), fact-checking URL recommendations (Vo & Lee, 2018), and reinforcement learning (Zhou & Zafarani, 2019).

2.2 Fake news detection based on user input

In addition to linguistic and visual qualities, user profile traits are employed as supplemental data to (Guo et al., 2018) identify fake news.

Study of implicit and explicit user profile feature has been used in this paper (Shu et al., 2019). Here author used two subsets, one is users spreading fake news and the other one is user spreading true news. Further he analyzed the relation of these user profiles. Feature that he has considered for his studies were: user verification, register time, political bias etc.

Another study (Zhou & Zafarani, 2019) looked into the social networks of users who disseminated news and illustrated their interrelationships. Fake news spreads faster than real news, spreaders connects with fake news more intensely than they do with legitimate news, and fake news spreaders form denser networks than real news spreaders, according to research that compared the frequency with which fake news stories spread and the total number of news stories spread by a user.

2.3 Fake News Detection Based on Propagation

Information propagation models make an effort to replicate the ways in which information spreads through time. (Vosoughi et al., 2017 and Shu et al., 2020) It has been shown that analysing the way news stories circulate on social media, such as through replies or retweets on Twitter, can help with the process of spotting fake information.

By constructing a news article propagation network, the authors proposed the HPFN model (Shu et al., 2020). The structural, temporal, and linguistic features were the three categories of features that the authors derived from the propagation network's two (Macro and Micro) levels. They looked at the hierarchical dissemination networks of fake and true news for structural, temporal, and linguistic characteristics.

Authors of the study (Silva et al., 2021) presented a technique that encodes the propagation tree, which allocates various degrees of priority to the nodes and cascades of the propagation tree. They also put forth a method for early fake news identification that reconstructs the useful information discovered in full-propagation trees made with early propagation trees.

In the following study, (Monti et al., 2019), textual node embedding features and graph were utilised to simulate the propagation network. The suggested method for predicting fake news consists of a softmax layer, two fully connected layers, two graph convolutional layers, and two fully connected layers.

In conclusion, the majority of the models that is now in use focus on a certain aspect of literature. In this work, we offer a method that dynamically uses propagation-based characteristics in addition to the static features as well implicit features that can identify whether given tweet is fake or real. The majority of the work that has been done focuses on identifying fake tweets, and experimentation is carried out precisely using the COVID-19 fake tweets as a case study.

3. A PRELIMINARY INVESTIGATION FOR DATA INSIGHT

A Preliminary Investigation of data has been performed in order to find patterns, identify anomalies, test hypotheses, and validate presumptions using summary statistics and graphical representations. To ascertain whether an author has sympathetic, hostile, or neutral attitude toward a particular topic, the technique of sentiment analysis involves computationally locating and classifying views in a text. In the first phase of analysis, polarity vs subjectivity among the Real and Fake tweet has been tested. The emotion that the sentence expresses is called polarity. It may also be neutral, good, or negative. The range of the float value for polarity is [-1, 1].Subjectivity is when text is an explanatory article which must be analysed in context. Figure 1(a) shows Polarity vs subjectivity in Real and Fake Tweets. The float value for polarity is in the range [-1,1], where 1 denotes a positive assertion and -1, a negative



one. Another float that falls between [0,1] is subjectivity.

Figure: 1(a) Polarity vs subjectivity in Real and Fake Tweets

Text Blob package for Python is used here for finding polarity and subjectivity. However no marginal difference found wrt polarity vs Subjectivity here in the fake and real tweets hence next experimentation is done on checking the Sentiments of Real and Fake tweet. Figure 1(b) shows here sentiments in real and fake content i.e. positive, negative or neutral sentiments involved in the sentence.

Observations for the Sentiments in Real and Fake Tweets:

No of Neutral > No of (Positive, Negative)

Neutral (fake) > Neutral (real)

Neutral (real, fake) > Positive (real, fake)

Neutral (real, fake) > Negative (real, fake)

Negative (fake) > Negative (real)

Positive (fake) > Positive (real)



Figure: 1(b) sentiments in real and fake tweet



Figure: 1(c) shows subjectivity score in Real and Fake Tweet

In the next experiment, subjectivity score of Real and Fake tweet has been observed as

shown in figure 1(c). Here Subjectivity score reflects that how many subjective or objective sentences are involved in Real and Fake Tweet. In the given figure, we can observe that the objectivity is more in case of fake tweets. Also the subjectivity is more in case of fake tweets. Objective statements pertain to factual information, but subjective sentences typically allude to personal opinion, emotion, or judgement. The float value of subjectivity is in the [0, 1] range. Figure 1(d) shows polarity and subjectivity score with flag. Here in figure 1(d), Value 0.7 indicates that there is more subjectivity, which eventually means that most of the material is public opinion rather than reality. Polarity is a float value that falls between [-1, 1], where 0.5 denotes a positive statement.

	polarity	subjectivity	sentiment_flag	subjectivity_flag
0	0.000000	0.000000	neutral	objective
1	0.000000	0.000000	neutral	objective
2	0.541667	0.708333	positive	subjective
3	-0.100000	0.200000	neutral	objective
4	0.250000	0.333333	neutral	neutral

Figure: 1(d) shows Polarity and subjectivity score with flag

3.1 The Experimental Dataset

The dataset used for the experimentation has taken from kaggle and the tweets are related to Covid-19. The dataset is available on the

https://www.kaggle.com/smid80/coronavirus-covid19-tweets-early-april. This dataset includes elements related to Twitter, including the tweet's title, the content of multiple tweets and the accounts that posted them, the hashtags used, and the accounts' geolocation. Retweets are not included in the dataset, but a variable that counts them is supplied. Along with the "retweet count," other features like "favorites count,""followers," and "friends" are also included in the dataset and have been effectively employed to increase the model's accuracy. Approximately 303692 tweets were used in the experimentation, of which 156612 were phoney and 147080 were actual.

	0	Unnamed: 0	313036 non-null	int64	
	1	status_id	313036 non-null	int64	
	2	user_id	313036 non-null	int64	
	3	created_at	313036 non-null	object	
	4	title	313036 non-null	object	textual
	5	text	313036 non-null	object	textual
ĺ	6	source	313035 non-null	object	
	7	is_quote	313036 non-null	bool	boolean
	8	is retweet	313036 non-null	bool	
	9	favourites_count	313036 non-null	int64	
	10	retweet_count	313036 non-null	int64	
	11	followers_count	313036 non-null	int64	numeric
	12	friends count	313036 non-null	int64	
ĺ	13	account created at	313036 non-null	object	
	14	label	313036 non-null	bool	

Figure 2: shows data types of dataset

So here Figure 2 shows the data types of the metadata. Here Title and Text are in object form which is a textual feature, is_quote, is_retweet are boolean features and

'friends_count', 'followers_count', 'retweet_count' and favourites_count' features are in numeric form.

4. METHODOLOGY

During our study, we have used different appraoches to handle the issue of Fake Tweet Detection. We have used tweet text data, user content features individually and in combination to find out the accuracy of detecting fake tweet. Alongwith this we have also tested derived attributes and calculated the conditional probability of this particular attribute indicating a real and fake news item

4.1 Text based Approach

Natural language processing (NLP) is the most obvious way to use pattern recognition to reliably identify fake news because it can extract information from the content of the tweet. The results of a model's construction depend on how effectively the data was preprocessed, which is a crucial stage. This method involves doing text normalisation, which involves changing all letters to lowercase or uppercase, turning numbers into words or eliminating numerals, removing punctuation, accent marks, and other diacritical marks, removing white spaces, extending abbreviations, and removing stop words. The frequency of a word in a document is referred to as term frequency. Words that appear too frequently across all papers are given lower ratings due to inverse document frequency. To put it simply, TF-IDF assigns a frequency score to words by emphasising those that occur more frequently within a single document, but not across several documents. The Tf-idf vectorizer reverses document frequency scores, tokenizes documents, learns vocabulary. Many deep learning algorithms employ the unsupervised learning method used Glove to generate word vector representations. To further assess the model's accuracy, basic classifier models including Naive Bayes, Logistic Regression, Decision Trees, Random Forests, and XG Boost and various deep learning models like recurrent neural networks and convolutional neural networks have been utilised.

4.2 Text and User content features based Approach

The pre-processed and normalised text used in the Tweet text feature is what allows the content's legitimacy to be checked. Through the use of NLP approaches, tweet text analysis has been possible to assess the tweet's reliability to some level. However, the text component might not be sufficient on its own to provide greater accuracy when it comes to determining believability. Quotes, favourites, retweets, followers, and friends are just a few examples of features that are included with a tweet and are referred to as its metadata. These features are important in the process of evaluating a tweet because they help it become more than just a string when certain metadata is added to it. As an

additional input dimension for an algorithm, this metadata is transformed into features.

In light of this, we have made an effort to evaluate the tweet's authenticity by looking at the media content, the account information, and the textual features. As a result, we have incorporated four user features—followers count,favorites count, friends count and retweet count along with tweet text processing in this instance, which can significantly alter the accuracy of the news prediction.

Here, we've attempted to develop a mechanism for evaluating the tweet's authenticity based on its correlation with the extra attributes. The classification of the tweet in this phase was carried out using a tweet and user content feature.

4.3 Probability vectors of derived features model

Text of the Tweets contains several such factors based on which we can predict the characteristics of tweet or we can predict whether the tweet is fake or real. In this method, we've added a heuristic approach to our original framework so that it can take the impact of the derived qualities into account. For data with attributes like username handles and URL domains, this method worked effectively. These characteristics enabled us to expand our present feature set with useful functionality. Here, attribute1 is the URL domain and attribute2 is the username handle. We estimate the conditional probability that this specific property distinguishes between real and fake news.

Domain Extraction:- Text of the Tweets contains several such factors based on which we can predict the characteristics of tweet or we can predict whether the tweet is fake or real. Domain Extraction does exactly the same task. In the Model, the urls have been extracted from text and calculated the fake and real probability of each url.

Fake and real probability of the url now became X factor to predict the fake or the real tweet. We then created a .json file to store them all. Figure 3(a) shows fake and real probability of extracted domains.

("https://t.co/oznS17Am8r": ("real probabil.	ity": 1.0,
"https://t.co/16MelKrx22": {"real probabili"	ty": 1.0,
"https://t.co/z5540%FS%D": {"real probabili;	ty": 1.0,
"https://t.co/vY4fVgAjuk": {"real probabili"	ty": 1.0,
"https://t.co/ozn517iLgT": ("real probabili"	ty": 1.0,
"https://t.co/iYrY5ciK15.": ("real probabil.	ity": 1.0,
"https://t.co/c18t005UV0": {"real probabili"	ty": 1.0,
"https://t.co/LZ3mOcDfNz": {"real probabili	ty": 1.0,
"https://t.co/pyHWGgFA4R": {"real probabili:	ty": 0.25,
"https://t.co/e6T5z11WYN": {"real probabili;	ty": 1.0,
"https://t.co/gumCU8RpTN": {"real probabili;	ty": 1.0,
"https://t.co/soPbCr9fIU": ("real probabili:	ty": 1.0,
"https://t.co/89qH2teUxe": {"real probabili:	ty": 1.0,
"https://t.co/d2xMleQU0y": ("real probabili:	ty": 1.0,
"https://t.co/PRvkFIgCVX": ("real probabili:	ty": 1.0,
"https://t.co/bvlyWzsbRh": ("real probabili"	ty": 1.0,
"https://t.co/BJszAcZgNG": {"real probabili;	ty": 0.333

Figure: 3(a) shows fake and real probability of extracted domains.

Username Extraction:- Username Extractiodeals with extracting unique usernames from the tweet text. We found the Fake and Real Probability of each username, which is again our factor on which our prediction relies. Then the creation of .json file is done to store them all together. Figure 3(b) shows fake and real probability of extracted

usernames.

'GovAbbott': {'fake_probability': 0.015844544095665172, 'real_probability': 0.9841554559043348, 'total_mentions': 3345}, 'GovernorTomWolf': {'fake_probability': 0.012422360248447204, 'real_probability': 0.9875776397515528, 'total_mentions': 1932}, 'GovParsonMO': {'fake_probability': 0.039544235924932974, 'real_probability': 0.960455764075067, 'total_mentions': 1492}, 'rondesantisfl': {'fake_probability': 0.0007987220447284345, 'real_probability': 0.9992012779552716, 'total_mentions': 1252}, 'govkemp': {'fake_probability': 0.002457002457002457, 'real_probability': 0.9975429975429976, 'total_mentions': 1221}, 'WHO': {'fake_probability': 0.5598678777869529, 'real_probability': 0.4401321222130471, 'total_mentions': 1211},

Figure: 3(b) shows fake and real probability of extracted usernames

By doing this, we are able to produce a probability vector that adds two more features to our new dataset. We gathered these characteristics from all of the tweets in our training set and determined how often each feature is real or fake based on the ground truth. As a result, for each new characteristic that are added to our existing feature set, we are able to construct a two-dimensional prediction vector

5. IMPLEMENTATION

5.1 Implementation details

Google colab8 is used to implement the hybrid deep learning model. Colab is a cloud environment for Jupyter notebooks that includes GPUs and TPUs for intensive computing. Python has been used to write the experiment code (pre-processing and classifiers). The hybrid CNN-RNN model has been implemented using the Keras Python10 package and the tensorflow module. Pandas library is used for reading the datasets whereas arrays are processed through Numpy library. For data pre-processing, the NLTK package is utilized. The Scikit-learn software is used to analyze the data, evaluate the results, and create baseline classifiers. Plotting graphs is done with the Matplotlib library.

5.2 Feature used for study

Feature	Name	Туре	Description			
TC1	Title	Textual	Highlight of the topic			
TC2	Text	Textual	Extended tweet section with details on the topic.			
TC3	Source	Textual	Shows the tweet's source.			
TC4	quote	Boolean	This setting decides whether a quote appears in the chosen tweet or not.			
TC5	retweet	Boolean	Retweeting sends a chosen tweet again. This field determines if the chosen tweet was re-tweeted.			

Tab.1 List of the Tweet content feature

Tab. 2 List of probability vectors

Feature	Feature Name	Feature type	Description
DC1	username_real	Numerical	Real username probability
DC2	username_fake	Numerical	Fake username probability
DC3	domain_real	Numerical	Real domain probability
DC4	domain_fake	Numerical	Fake domain probability

Tab. 3 List of User content feature

Feature	Name	Туре	Description
UC1	favourites_ count	Numerical	It refers to how many tweets a single user has marked as favourites or how many tweets this user has liked overall over the account's existence.
UC2	Retweet_ Count	Numerical	Only the original tweet is counted when re- tweeted.
UC3	Followers_ Count	Numerical	Number of the Twitter account's followers.
UC4	Friends_ Count	Numerical	The total number of Twitter friends is shown. But it might be possible to learn something about how the twitter account is being used from the ratio of followers to friends.



Figure 4: Block diagram of CRED_Tweet model using Tweet Text, Content & Derived features

5.3 CRED Tweet Implementation using Tweet Content & Derived features

In this approach, we handled multiple inputs simultaneously which have textual data (tweet text), continuous data i.e metadata in numeric form and derived inputs which consists of username and domain probability. Figure 4 shows block diagram of CRED_Tweet model using Tweet Content & Derived features.

When we look at the internal structure of the model, tweet text input, user content and derived input has been passed through a separate neural network. In the first phase, text input is given to the Data Preprocessing layer. The various text preprocessing steps such as Tokenization, Lower casing, Stop words removal, Stemming/Lemmatization. The text analysis approach based on Convolutional neural network (CNN) can collect significant text properties by pooling, whereas the LSTM model predicts well based on contextual information. This text input has been processed through the convolution kernel and concatenated at the end.

In the next phase, user contents has been taken from tweet metadata. Tweet metadata such as friends count, followers count, favourites count and retweet count has been processed through separate process. In this approach, we handle multiple inputs simultaneously which have textual data (tweet text), continuous data i.e metadata in numeric form.

Lastly we have used derived attributes like URL domains, username handles. We have extracted these attributes from the tweet text and calculated the conditional probability of this particular attributes and used as an additional parameter along with Tweet text and content features. Fake and real probability of these attributes becomes X factor here to predict the fake or the real tweet.

These Text contents, User contents and probability vectors has been processed through separate CNN layers and concatenated with processed text input vectors. A fully connected layer has been used before the prediction is produced. Lastly the output layer

i.e. dense layer is utilized in conjunction with a single neuron and a softmax activation function to predict whether a tweet is fake or real. Since this is a binary classification issue, binary_crossentropy has been used as a loss function and ADAM adaptive optimizer algorithm is used inorder to improve the results. As compared to other models, CRED_tweet model gives greater test accuracy with approximately lesser training time and same weights.

5.4 Enhanced Algorithm for CRED_Tweet

Approach

Phase I: Extracting text features

In the first phase, text input is given to the word embedding layer followed by a convolution kernel. Here the Convolution kernel consists of one dimensional convolution, activation and the consolidated features are then sent to the LSTM through the max pooling layer. The text analysis approach based on Convolutional neural network (CNN) can collect significant text properties by pooling, whereas the LSTM model predicts well based on contextual information.

Assuming that the maximum length of the tweet is n, let Z be the text input di TC (text features).

$$Z_{i}^{TC} = Z_{1} \oplus Z_{2} \oplus Z_{3} \oplus \dots Z_{n}$$
 (6)

Equation (7) displays the default convolution function.

$$Z = W^{T}.X + b$$
 (7)
$$Z := f(W^{T}.Z_{i}^{TC} + b)$$
 (8)

The feature Zi presented in equation (8 is a result of the word matrix processing through a convolution layer. The bias is represented by the letter b in this instance.

$$Z_{i} = \max\{0, Z_{i}\}$$
 (9)

Zi has been passed through three distinct convolution layers made up of 128 neurons during the experiment. To utilise all of the feature map's potential. Zi' represents the highest value.

 $Z_{i}(ht) = Z_{i}(hft \oplus hbt)$ (10)

These token vectors (Z i) are further encoded using a Bi-LSTM, using the forward and backward layers which processes the N vectors in opposite directions. a hidden state hft is emitted by the forward LSTM at each time-step, which is concatenated with the corresponding hidden state hbt of the backward LSTM

Phase II: Extracting derived features adding as additional metadata

 $D_i^{DC} = D_1 \oplus D_2 \oplus D_3 \oplus D_n$ (11)

Let D be the probability input \in di DC. Here 4 derived features has been used (real username, fake username, real domain, and fake domain)

$$D_{i}=f(W^{T}. D_{i}^{DC} + b)$$
 (12)

 $D_{i} = max\{0, D_{i}\}$ (13)

Phase III: Extracting metadata features, concatenating with text and derived metadata

$$X_{i}^{UC} = X_{1} \oplus X_{2} \oplus X_{3} \oplus \dots X_{n} \quad (14)$$

Let X be the metadata numeric input \in di UC. Xi UC has been passed through 2 different convolution layers of 128 neurons

$$X_{i}=f(W^{T}, X_{i}^{UC} + b)$$
 (15)
 $X_{i}' = max\{0, X_{i}\}$ (16)

Concatenate all the features (text, derived and metadata) feature

 $Z_{i}^{'} = Z_{i(ht)} \bigoplus X_{i}^{'} \bigoplus D_{i}^{'}$ (17)

We predict for y for given di TC, di DC, di UC ; where θ represents the parameters of the model used during the experimentation time. activation function(f) on (Zi')

y_i=f(Zi') (18)

Pred (yi for given di TC, di UC, di DC; θ) = activation function(f) on (Zi')

6. EVALUATION AND RESULTS

The purpose of the methodology is to see how well our various feature-based techniques distinguish between fake and real tweets. In our first experiment, we have used exclusively Text features based approach and performed experimentation using different machine learning algorithm i.e. Random Forest, Logistic Regression, Decision Tree, Naïve Bayes and XG-Boost algorithm with tf-idfvectorizer for tweet features of the tweet and deep learning algorithm such as BiLSTM, CNN and hybrid model was used. The Evaluation parameters used here were accuracy, precision, recall (sensitivity). Table 5 shows the results of ML Classifier and Table 6 shows deep neural network results with Text feature. Among all the algorithms CNN-LSTM model outperformed in terms of accuracy.

Algorithm	Precision	Recall	Accuracy
Logistic Regression using count vectorizer	0.92	0.92	0.924
Logistic Regression with TF-IDF	0.93	0.93	0.925
Random forest with Count vectorizer	0.92	0.93	0.926
Random forest with tf_idf	0.92	0.93	0.928
Naïve bayes using count vectorizer	0.92	0.92	0.922
Naïve bayes using tf-idf	0.93	0.92	0.923

Tab. 5 Experimental Results with ML Classifier

Algorithm	Precision	Recall	Accuracy
Bidirectional LSTM	0.89	0.90	0.896
Simple CNN Model	0.92	0.92	0.929
Modified CNN Model	0.93	0.93	0.934
CNN-LSTM Model	0.94	0.94	0.942

 Tab. 6: Experimental Results with Neural Network

Table 7 shows the performance results of various metadata features. Features used for the experimentation were is quote, is retweet, favourites count, followers count, friends count and retweet count. We have also tried to use different combination of these features inorder to check interrelation between these feature and in terms of achieving higher accuracy among all. We have used CNN model to evaluate the metadata features based on different evaluation parameters such as accuracy, F1 score, precision and recall.

Feature used	Feature used	Accuracy	F1 Score	Precision	Recall
favourites_count	UCF1	92.47	71.25	87.52	30.12
retweet_count	UCF2	92.46	69.26	84.05	28.49
followers_count	UCF3	93.48	73.65	97.55	34.34
friends_count	UCF4	92.43	67.34	93.97	24.39
favourites_count','retweet_co unt	UCF1, UCF2	92.5	68.01	92.67	25.53
retweet_count, 'followers_count	UCF2, UCF3	92.45	67.96	93.88	24.59
followers_count','friends_coun t	UCF3, UCF4	93.45	73.59	96.7	34.35
'is_quote','is_retweet'	TCF4, TCF5	92.49	67.52	95.62	24.53
'is_quote','is_retweet','favourit es_count','retweet_count','foll owers_count','friends_count'	TCF4, TCF5, UCF1, UCF2, UCF3, UCF4	93.22	75.34	79.96	41.15
'favourites_count','retweet_co unt','followers_count','friends_ count'	UCF1, UCF2, UCF3, UCF4	94.87	85.23	74.92	71.74

 Tab 7 : shows accuracy parameters for metadata features

Further to this, alongwith the various evaluation parameters, accuracy, precision, recall (sensitivity), F1-score we have used true negative rate (specificity), true positive rate (TPR), PRC (precision-recall curve), ROC (receiver operating curve), and other metrics to assess the projected results. We've looked at how various user profile categories, tweet content components, and a mix of both performed. Here Table 8 shows accuracy parameters for Machine Learning classifiers with text and metadata features in standalone mode as well as in the combination. Table 9 shows accuracy parameters for CRED_Tweet Approach with different features such as CRED_Tweet Approach with Tweet Text Content (TC), combining Tweet Text Content (TC) and User profile Content (UC) and lastly combination of Tweet Text Content (TC), User profile Content (UC), Derived content (DC) which outperforms among all the models and gives the 99.42 % F1 score.

Approaches	Measures	LR	NB	RF	DT	XG Boost
	Precision	90.68	98.91	91.0	79.62	96.37
	Recall	75.07	57.68	74.4	76.48	64.54
Text Based Approach	F1 score	83.39	77.37	83.21	78.19	80.35
Tweet Text Content (TC)	ROC	0.91	0.88	0.91	0.78	0.89
	PRC	0.93	0.91	0.93	0.73	0.91
	Accuracy	83.48	78.26	83.32	78.19	80.83
	Precision	86.23	91.70	93.78	83.79	82.47
	Recall	34.81	72.10	82.56	85.02	67.3
Metadata Based Approach	F1 score	72.96	82.15	93.3	84.40	85.81
User profile Content (UC)	ROC	0.94	0.81	0.99	0.92	0.97
	PRC	0.77	0.79	0.94	0.73	0.85
	Accuracy	90.17	91.80	90.79	91.96	91.46
	Precision	85.17	91.80	91.72	92.13	90.91
	Recall	91.02	74.40	90.97	92.54	91.41
Combining Text and Metadata of the	F1 score	87.98	84.15	91.64	92.55	91.41
profile Content (LC)	ROC	0.93	0.83	0.97	0.93	0.97
	PRC	0.94	0.81	0.97	0.89	0.97
	Accuracy	87.98	84.39	91.65	92.56	91.42

 Tab: 8 shows accuracy parameters for Machine Learning classifiers with text and metadata features

Measures	CRED_Tweet Approach with Tweet Text Content (TC)	CRED_Tweet Approach with Tweet Text Content (TC) + User profile Content (UC)	CRED_Tweet Approach with Tweet Text Content (TC) + User profile Content (UC)+Derived content (DC)
Precision	94.00	98.44	99.78
Recall	94.00	96.56	99.03
F1 score	94.20	97.60	99.42
Accuracy	94.20	97.60	99.42

Tab: 9 shows accuracy parameters for CRED_Tweet Approach with different features

7. CONCLUSION

Our proposed approach named CRED_Tweet comprises of 3 components i.e tweet text analysis, metadata analysis i.e user content and derived features. We investigated text features of the tweet using various natural language processing methods. During metadata analysis, we studied the inter relation among the different user content features and used them efficiently in our algorithm. The probability vector analysis is done to calculate the probability of username and domains included in tweets. Finally, we efficiently aggregated all the components successfully to determine the veracity of the input tweet. This work can be naturally extended by classifying tweets based on their subjects and examining the efficacy of extracted features for fake tweet identification on each topic independently using our suggested method. In order to detect fake tweet, the use of tweet text, metadata and derived features with deep neural network seems promising.

Compared to other industry-standard deep learning and machine learning methods, our suggested model performs better. In the future, we'll keep examining neural network topologies that are more intricate than CNN and Bi-LSTM. They can be incredibly helpful when used in conjunction with task-specific function engineering strategies. Even though there are a tons of studies on fake news identification and detection already, there is always potential for improvement. New, fundamental insights into the nature of fake news can result in models that are more accurate and effective.

We hereby declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

 Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake news detection on social media: A data mining perspective. ACM SIGKDD Explorations Newsletter, 19(1), 22–36.

https://doi.org/10.1145/3137597.3137600

[2] Vishwakarma, D. K., & Jain, C. (2020). Recent state-of-the-art of fake news detection: A review. 2020 International Conference for Emerging Technology, INCET 2020, 1–6

https://doi.org/10.1109/INCET49848.2020.9153985.

- [3] Zhang, X., & Ghorbani, A. A. (2019). An overview of online fake news: Characterization, detection, and discussion. Information Processing and Management. https://doi.org/10.1016/j.ipm.2019.03.004"
- [4] Bondielli, A., & Marcelloni, F. (2019). A survey on fake news and rumour detection techniques. Information Sciences, 497, 38–55. https://doi.org/10.1016/j.ins.2019.05.035
- [5] Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., Zittrain, J. L. (2018). The science of fake news. Science, 359(6380), 1094– 1096.https://doi.org/10.1126/science.aao2998
- [6] Varshney, D., & Vishwakarma, D. K. (2021a). A review on rumour prediction and veracity assessment in online social network. Expert Systems with Applications, 168, Article 114208. https://doi.org/10.1016/j.eswa.2020.114208
- [7] Meel, P., & Vishwakarma, D. K. (2020). Fake news, rumor, information pollution in social media and web: A contemporary survey of state-of-the-arts, challenges and opportunities. Expert Systems with Applications, 153, Article 112986. https://doi.org/10.1016/j.eswa.2019.112986
- [8] Rubin, V. L., Chen, Y., & Conroy, N. J. (2015). Deception detection for news: Three types of fakes. Proceedings of the Association for Information Science and Technology, 52(1), 1–4. https://doi.org/10.1002/pra2.2015.145052010083
- [9] Gautam Kishore Shahi, Anne Dirkson, Tim A. Majchrzak, An exploratory study of COVID-19 misinformation on Twitter, Online Social Networks and Media, Volume 22, 2021, 100104, ISSN 2468-6964, https://doi.org/10.1016/j.osnem.2020.100104.
- [10] M. Cinelli, W. Quattrociocchi, A. Galeazzi, C.M. Valensise, E. Brugnoli, A.L.Schmidt, P. Zola, F. Zollo, A. Scala, The COVID-19 social media infodemic, arXiv(2020 (Preprint)).
- [11] R. Kouzy, J. Abi Jaoude, A. Kraitem, M.B. El Alam, B. Karam, E. Adib, J. Zarka, C. Traboulsi, E. Akl, K. Baddour, Coronavirus goes viral: quantifying the COVID- 19 misinformation epidemic on Twitter, Cureus 12 (2020) e7255, https://doi.org/10.7759/cureus.7255.
- [12] R. Gallotti, F. Valle, N. Castaldo, P. Sacco, M. De Domenico, Assessing the risks of "infodemics" in response to COVID-19 epidemics, preprint (2020).
- [13] L. Singh, S. Bansal, L. Bode, C. Budak, G. Chi, K. Kawintiranon, C. Padden, R.Vanarsdall, E. Vraga, Y. Wang, A first look at COVID-19 information and misinformation sharing on Twitter, preprint (2020).
- [14] K.-C. Yang, C. Torres-Lugo, F. Menczer, Prevalence of low-credibility information on twitter during the COVID-19 outbreak, ArXiv preprint (2020). 10.36190/2020.16.
- [15] J. Allen, B. Howland, M. Mobius, D. Rothschild, D.J. Watts, Evaluating the fakenews problem at the scale of the information ecosystem, Sci. Adv. 6 (14) (2020) eaay3539, https://doi.org/10.1126/sciadv.aay3539.

- [16] F. Zollo, A. Bessi, M. Del Vicario, A. Scala, G. Caldarelli, L. Shekhtman, S. Havlin, W. Quattrociocchi, Debunking in a world of tribes, PLoS One 12 (7) (2017) e0181821, https://doi.org/10.1371/journal.pone.0181821.
- [17] A. Bovet, H.A. Makse, Influence of fake news in Twitter during the 2016 US presidential election, Nat. Commun. 10 (1) (2019) 1–14, https://doi.org/10.1038/s41467-018-07761-2.
- [18] N. Grinberg, K. Joseph, L. Friedland, B. Swire-Thompson, D. Lazer, Politicalscience: fake news on Twitter during the 2016 U.S. presidential election, Science363 (6425) (2019) 374–378, https://doi.org/10.1126/science.aau2706.
- [19] C. Shao, P.-M. Hui, L. Wang, X. Jiang, A. Flammini, F. Menczer, G. Luca Ciampaglia, Anatomy of an online misinformation network, PLoS One 13 (4) (2018) e0196087, https://doi.org/10.1371/journal.pone.0196087.
- [20] Zhang, X., & Ghorbani, A. A. (2019). An overview of online fake news: Characterization, detection, and discussion. Information Processing and Management. https://doi.org/10.1016/j.ipm.2019.03.004"
- [21] Oshikawa, R., Qian, J., & Wang, W. Y. (2018). A Survey on Natural Language Processing for Fake News Detection. http://arxiv.org/abs/1811.00770.
- [22] Zhou, & Zafarani, R. (2019). Network-based Fake news detection: A pattern-driven approach. ACM SIGKDD Explorations Newsletter, 21(2), 48–60. https://doi.org/10.1145/3373464.3373473
- [23] Ghosh, S., & Shah, C. (2019). Toward Automatic Fake News Classification. Proceedings of the 52nd Hawaii International Conference on System Sciences, 6, 2254–2263. https://doi.org/10.24251/hicss.2019.273.
- [24] Choudhary, A., & Arora, A. (2021). Linguistic feature based learning model for fake news detection and classification. Expert Systems with Applications, 169, Article 114171.https://doi.org/10.1016/j.eswa.2020.114171
- [25] Potthast, M., Kiesel, J., Reinartz, K., Bevendorff, J., & Stein, B. (2018). A stylometric inquiry into hyperpartisan and fake news. ACL 2018–56th Annual Meeting of the Association for Computational Linguistics. Proceedings of the Conference (Long Papers), 1, 231–240. https://doi.org/10.18653/v1/p18-1022.
- [26] Zhang, X., & Ghorbani, A. A. (2019). An overview of online fake news: Characterization, detection, and discussion. Information Processing and Management. https://doi.org/10.1016/j.ipm.2019.03.004
- [27] Dungs, S., Aker, A., Fuhr, N., & Bontcheva, K. (2018). Can Rumour Stance Alone Predict Veracity?. In Proceedings of the 27th International Conference on Computational Linguistics (pp. 3360–3370).
- [28] Qian, F., Gong, C., Sharma, K., & Liu, Y. (2018). Neural user response generator: Fake news detection with collective user intelligence. IJCAI International Joint Conference on Artificial Intelligence, 2018-July, 3834–3840. https://doi.org/10.24963/ijcai.2018/533.
- [29] Hosseinimotlagh, S., & Papalexakis, E. E. (2018). Unsupervised Content-Based Identification of Fake News Articles with Tensor Decomposition Ensembles. Cs: Ucr.Edu. https://doi.org/10.475/123.
- [30] Kaliyar, R. K., Goswami, A., & Narang, P. (2021). DeepFakE: Improving fake news detection using tensor decomposition-based deep neural network. Journal of Supercomputing, 77(2), 1015–1037. https://doi.org/10.1007/s11227-020-03294-y
- [31] Gupta, S., Thirukovalluru, R., Sinha, M., & Mannarswamy, S. (2018). CIMT Detect: A community infused matrix-tensor coupled factorization based method for fake news detection. In Proceedings of

the 2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining. https://doi.org/10.1109/ASONAM.2018.8508408

- [32] Long, Y., Lu, Q., Xiang, R., Li, M., & Huang, C.-R. (2017). Fake news detection through multiperspective speaker profiles. Proceedings of the Eighth International Joint Conference on Natural Language Processing, 2(8), 252–256. http://www.aclweb.org/anthology/117-2043.
- [33] Song, C., Ning, N., Zhang, Y., & Wu, B. (2021). A multimodal fake news detection model based on crossmodal attention residual and multichannel convolutional neural networks. Information Processing and Management, 58(1), Article 102437. https://doi.org/10.1016/j.ipm.2020.102437
- [34] Girgis, S., & Gadallah, M. (2018). Deep Learning Algorithms for Detecting Fake News in Online Text. In 13th International Conference on Computer Engineering and Systems (ICCES) (pp. 93–97). https://doi.org/10.1109/ICCES.2018.8639198
- [35] Liang, G., He, W., Xu, C., Chen, L., & Zeng, J. (2015). Rumor identification in microblogging systems based on users' behavior. IEEE Transactions on Computational Social Systems, 2(3), 99–108. https://doi.org/10.1109/TCSS.2016.2517458
- [36] Chen, Y., Hu, L., Sui, J., & Gong, W. (2019). Attention-residual network with CNN for rumor detection. International Conference on Information and Knowledge Management, Proceedings, 1121–1130. https://doi.org/10.1145/3357384.3357950
- [37] Meel, P., & Vishwakarma, D. K. (2021a). A temporal ensembling based semi-supervised ConvNet for the detection of fake news articles. Expert Systems with Applications, 177(April), Article 115002. https://doi.org/10.1016/j.eswa.2021.115002
- [38] Wang, Ma, F., Jin, Z., Yuan, Y., Xun, G., Jha, K., Su, L., & Gao, J. (2018). EANN: Event adversarial neural networks for multi-modal fake news detection. Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,849–857. https://doi.org/10.1145/3219819.3219903.
- [39] Vo. N., & Lee, K. (2018). The rise of guardians: Fact-checking URL recommendation to combat fake news. In 41st International ACM SIGIR Conference on Research and Development in Information Retrieval. https://doi.org/10.1145/3209978.3210037
- [40] Chen, Y., Hu, L., Sui, J., & Gong, W. (2019). Attention-residual network with CNN for rumor detection. International Conference on Information and Knowledge Management, Proceedings, 1121–1130. https://doi.org/10.1145/3357384.3357950
- [41] Guo, H., Cao, J., Zhang, Y., Guo, J., & Li, J. (2018). Rumor detection with hierarchical social attention network. International Conference on Information and Knowledge Management, Proceedings, 943– 952. https://doi.org/10.1145/3269206.3271709
- [42] Shu, K., Zhou, X., Wang, S., Zafarani, R., & Liu, H. (2019). The Role of User Profile for Fake News Detection. ASONAM '19. In Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (pp.436–439). https://doi.org/10.1145/3341161.3342927
- [43] Vosoughi, S., Mohsenvand, M. N. E. O., & Roy, D. E. B. (2017). Rumor gauge : predicting the veracity of rumors on Twitter r r. ACM Transactions on Knowledge Discovery from Data, 11(4). https://doi.org/10.1145/3070644
- [44] Shu, K., Mahudeswaran, D., Wang, S., Lee, D., & Liu, H. (2020). FakeNewsNet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. Big Data, 8(3), 171–188. https://doi.org/10.1089/big.2020.0062

- [45] Silva, A., Han, Y., Luo, L., Karunasekera, S., & Leckie, C. (2021). Propagation2Vec: Embedding partial propagation networks for explainable fake news early detection. Information Processing and Management, 58(5), Article 102618. https://doi.org/10.1016/j.ipm.2021.102618
- [46] Monti, F., Frasca, F., Eynard, D., Mannion, D., & Bronstein, M. M. (2019). Fake news detection on social media using geometric deep learning. https://arxiv.org/abs/1902.06673.