



ISSN:2147-6799

www.ijisae.org

Original Research Paper

Combating Fake News on Twitter: A Machine Learning Approach for Detection and Classification of Fake Tweets

Archana Nanade^{1*}, Dr. Arun Kumar²

Submitted: 25/08/2023 Revised: 11/10/2023 Accepted: 26/10/2023

Abstract: The proliferation of false information across social media has emerged as a significant challenge in recent times. Within this landscape, the platform formerly known as Twitter, now referred to as X, stands as a widely used medium for disseminating news and data, rendering it especially susceptible to the proliferation of misleading or fabricated content. In response to this issue, this research paper presents an innovative strategy hinging on BERT (Bidirectional Encoder Representations from Transformers) technology to address the propagation of deceptive news on Twitter. By means of meticulously acquired Twitter data pertaining to news items, the dataset undergoes careful pre-processing to render it compatible with the BERT model. This data is subsequently partitioned into training and validation subsets. The BERT model is then meticulously fine-tuned via a feed-forward neural network and optimized leveraging the Adam optimizer. Throughout the training process, vigilance is maintained over loss values, augmented by techniques such as dropout and regularization to enhance generalization capabilities. The final model selection is dictated by the validation loss metric. Harnessing the capabilities of advanced natural language processing methods, this endeavour contributes to the evolution of robust instruments for unearthing misleading content, fostering a more dependable informational milieu across social media domains with an accuracy of 77.29%. The study's conclusions carry implications for elevating the integrity and credibility of news propagation within the digital epoch

Keywords: BERT, Twitter, Fakes News detection

1. Introduction

The rise of social media has transformed information sharing, offering instant global connectivity. However, it has also fuelled misinformation, particularly on platforms like Twitter due to real-time sharing and a wide user base. Twitter is a prominent platform for news and opinions due to its concise format and rapid reach. Yet, this also makes it susceptible to fake news due to limited context and swift amplification. Detecting and countering Twitter fake news is vital. False information risks misleading the public, distorting discourse, and eroding trust. Ensuring accurate information is crucial for informed choices, democratic processes, and news credibility.

The proliferation of false information across social media has emerged as a significant challenge in recent times. Research studies by (1) and (2) have revealed the rapid spread of fake news on platforms like Twitter, leading to widespread misinformation and confusion. The real-time nature of Twitter, characterized by short messages (tweets) and quick sharing, makes it particularly prone to the dissemination of misleading content.

This paper confronts the urgent challenge of recognizing fraudulent news on Twitter by introducing a pioneering algorithmic approach. The primary objective of this paper

 ¹ Sir Padampat Singhania University, Udaipur, Rajasthan-ORCID ID: 0000-0002-2886-5810
 ² Sir Padampat Singhania University, Udaipur, Rajasthan ORCID ID: 0000-0003-1972-5717
 * Corresponding Author Email: archana.nanade@spsu.ac.in focuses on devising a resilient and effective system capable of proficiently pinpointing and categorizing deceptive content, thereby contributing to the broader endeavour of cultivating a more dependable informational environment on the Twitter platform.

To tackle this problem, this paper proposes an algorithm based on the principles of natural language processing and deep learning. The proposed approach leverages the power of the BERT (Bidirectional Encoder Representations from Transformers) model (3), a state-of-the-art language representation model known for its ability to capture contextual information effectively. By harnessing the advanced capabilities of the BERT model, the aim is to enhance the accuracy and efficiency of fake news detection on Twitter.

The algorithm proposed in this paper follows a systematic pipeline tailored to the characteristics of Twitter data. The approach begins by collecting a comprehensive dataset of tweets containing news-related content. This dataset serves as the foundation for training and evaluating the model. To prepare the data, tokenization and normalization techniques are employed, converting the raw text into a suitable format for analysis and classification

Next, the dataset is split into training and validation sets, ensuring the model's ability to generalize beyond the specific instances encountered during training. The BERT model is then trained using the training data, fine-tuning it specifically for fake news detection. During the training process, optimization algorithms are employed, and hyperparameters are carefully selected to optimize the model's performance.

After training, the effectiveness of the proposed approach is evaluated using robust metrics such as accuracy, precision, recall, and F1 score. By comparing the results with existing methods and benchmark datasets, the paper aims to validate the algorithm's performance and demonstrate its superiority in detecting fake news on Twitter.

The findings from this research provide valuable insights into the challenges of combating fake news on Twitter and offer practical solutions for promoting the authenticity and reliability of information. By leveraging advanced techniques and the power of deep learning, this paper contributes to the ongoing efforts in addressing the global issue of fake news.

This paper presents an innovative algorithm for detecting fake news on Twitter using the BERT model. By harnessing the capabilities of natural language processing and deep learning, the proposed approach aims to enhance the accuracy and efficiency of fake news detection, ultimately contributing to the creation of a more trustworthy information ecosystem on Twitter. The subsequent section 3 of this paper delves into the methodology, experimental setup, results, and discussions, providing a comprehensive analysis of the approach's performance and its implications for combatting fake news on social media platforms, with a specific focus on Twitter.

2. Literature Survey

The detection and mitigation of fake news on social media platforms, particularly Twitter, have garnered significant attention from researchers and practitioners in recent years. Several studies have explored various approaches and techniques to address this critical issue. In this literature survey, we present a summary of the existing research related to fake news detection on Twitter, focusing on the use of advanced algorithms and models.

The widespread dissemination of fake news on social media platforms, particularly on Twitter, has raised concerns about the credibility and reliability of the information shared. Researchers have been actively working on developing effective approaches to detect and combat fake news in order to promote accurate information and informed decision-making.

(4) proposed FakeBERT, a deep learning approach based on BERT, a state-of-the-art language model, for fake news detection on Twitter. By training the model to learn the linguistic features specific to fake news, they achieved a high accuracy of 98.90%. This highlights the potential of BERT-based models in capturing the nuanced language patterns associated with fake news.

In the context of the COVID-19 pandemic, (5) developed a BERT-based model capable of classifying COVID-19 related tweets into fake news and real news categories. They also incorporated sentiment analysis to gauge the emotional tone of the tweets. Their model achieved an accuracy of 90% for fake news detection and 93% for sentiment analysis, showcasing its effectiveness in tackling fake news surrounding a specific topic.

(6) proposed a fusion technique-based ensemble deep learning model that combined two powerful models, CT-BERT and RoBERTa, to detect fraudulent tweets related to the COVID-19 epidemic. With an accuracy of 98.88%, their model demonstrated the potential of leveraging multiple deep learning models for enhanced fake news detection. (7) introduced BERTweet, a pre-trained language model specifically designed for English tweets. BERTweet effectively captures the linguistic features unique to tweets, enabling it to be a valuable tool for various natural language processing tasks, including fake news detection.

(8) focused on a comprehensive NLP-based approach for fake news detection on Twitter. By employing techniques such as named entity recognition, sentiment analysis, and hate speech detection, they developed a robust model capable of identifying fake news tweets based on multiple linguistic cues.

Other studies, such as (9) and (10), and (11), explored machine learning and feature engineering techniques for fake news detection on Twitter. These approaches incorporated factors like user engagement metrics (e.g., retweets, likes) and linguistic features to distinguish between genuine and fake news tweets.

(12) explored the importance of neural network initialization and its impact on the expressiveness and performance of deep learning models. Their findings shed light on the nuances of initializing models effectively, which is crucial for improving the accuracy and generalizability of models used for fake news detection. (13) introduced the use of recurrent neural networks (RNNs) for detecting rumors in microblogs. Their work underscores the potential of sequence modeling techniques to capture temporal dependencies in tweets, enabling more accurate identification of fake news by considering the context in which they are shared.

(14) delved into the psychological factors that influence individuals' susceptibility to fake news. By analyzing cognitive processes, political beliefs, and other psychological dimensions, their study contributes to understanding the cognitive biases that underlie the spread and acceptance of fake news on social media platforms like Twitter. (15) examined the credibility of information shared on Twitter. Their work emphasized the significance of social network analysis and user behavior in assessing the credibility of tweets, which can complement the technical approaches for fake news detection and provide a more holistic understanding of the problem.

(16) presented a hybrid deep learning model for fake news detection. Their approach combines convolutional neural networks (CNNs) and long short-term memory (LSTM) networks, demonstrating the potential of leveraging multiple neural architectures to enhance the accuracy of fake news detection. The research by (17) focused on fake news detection on Twitter and presents an approach that combines feature engineering and machine learning techniques. Their work contributes to the early exploration of machine learning methods for identifying fake news on social media platforms like Twitter. By implementing feature engineering, the authors demonstrate their commitment to crafting effective feature sets for machine learning models, which is a crucial aspect of fake news detection. This study serves as a foundational reference for understanding the early developments in the field.

The research conducted by (18) focused on detecting online fake news using N-gram analysis and machine learning techniques. N-grams are an important text analysis tool, and their utilization in fake news detection underscores the relevance of linguistic features. This study showcases the application of machine learning in distinguishing fake news from legitimate content, further emphasizing the importance of computational approaches in the fight against misinformation. (19) presented a hybrid approach for detecting spammers in Twitter marketing, which has relevance in the broader context of content credibility. The use of social media analytics and bio-inspired computing signifies the multifaceted nature of combating misinformation. By considering marketing-related content and incorporating bio-inspired computing techniques, this research provides valuable insights into addressing malicious content and its impact on online platforms.

(20) delved into the significance of social context in fake news detection. It highlights that detecting misinformation goes beyond analyzing the content alone and involves understanding the broader social context in which information is shared. This study underscores the importance of considering factors like user behavior and interaction patterns in identifying fake news, offering a holistic perspective on the problem.

(21) research focused on text classification and sentiment prediction on social media datasets using a multichannel convolutional neural network (CNN). While not directly related to fake news detection, this reference is relevant because it demonstrates the utility of deep learning techniques for processing social media text data. Such techniques can potentially be applied in fake news detection systems to enhance their performance. Gupta and Kumaraguru's research [22] focused on the credibility ranking of tweets during high-impact events. This study highlights the importance of assessing the credibility of information on social media platforms, particularly during critical events when misinformation can have significant consequences. By addressing credibility ranking, the paper contributes to the broader understanding of information quality and trustworthiness on Twitter, which is closely related to fake news detection. (23) explored user perception of information credibility of news on Twitter. This work investigates how users perceive and evaluate the credibility of news shared on the platform. Understanding user perceptions is crucial in fake news detection, as it sheds light on the human side of the problem. Recognizing that people make judgments about information credibility informs the design of more effective detection algorithms and strategies. (24)'s research addressed the challenge of finding true and credible information on Twitter. This work underscores the importance of distinguishing reliable information from the vast amount of data circulating on the platform. The research likely includes valuable insights and techniques for assessing credibility that can be adapted for fake news detection, as both areas share the goal of differentiating between trustworthy and untrustworthy content.

The existing literature offers valuable insights into the detection of fake news on Twitter. While various approaches and models have been explored, our research paper contributes by proposing an innovative algorithm based on the BERT model, specifically tailored for fake news detection on Twitter. By building upon the advancements in NLP and deep learning, our approach aims to enhance the accuracy and efficiency of detecting and classifying fake news, ultimately contributing to the creation of a more trustworthy information ecosystem on Twitter.

3. Methodology and Algorithm

The methodology of this research paper follows a four-part approach, as shown in Figure 1, for detecting fake news using a modified BERT based algorithm. This modified algorithm is named as TweetTruth. In the first part, explained in section 3.1, the collected data is subjected to thorough cleaning and preprocessing, including steps such as noise removal, handling missing values, and text normalization. Tokenization is then applied to convert the pre-processed text into tokens suitable for further analysis. The second part focuses on creating a data loader that efficiently handles the tokenized data, including batching, shuffling, and loading onto the appropriate hardware, as explained in section 3.2. The third part involves creating a BERT model specifically for fake news detection, finetuning it using the appropriate architecture and hyperparameters, as explained in section 3.3. Lastly, as

explained in section, 3.4, the performance of the trained model is evaluated using various metrics on a validation or test dataset, assessing its ability to accurately classify fake news



Fig 1. Flow of TweetTruth Model

3.1 The Data Cleaning and Preprocessing for Tokenization

This section discusses the pseudo code steps used in TweetTruth model.

Data Gathering: It is very crucial that the collected dataset be accurate and diverse. For the acquisition of datasets,(25), is used as it has an enormous collection of various datasets. To read the dataset and pre-process it for further operations Pandas - Python package (26) is used for data analysis and associated manipulation of tabular data in data frames.

Furthermore, the dataset is divided into two classes of samples to achieve a diversity of tweet types. One with no real disaster and the other with a real disaster. This will be useful for segregation as well as training. The acquired and pre-processed dataset will be sent for further clarification. Table 1 lists the functions used in the algorithm and a short description for each

Table 1.	Functions	for	algorithm	for	pre-processing
					r · r · · · · · · · · · · · · · · · · ·

Function	Description
Load(x)	load x into to system
Clean $(x, y_1, y_2 y_i)$	remove y_1 , y_2 y_i attribute from x
Split(x, set1, set2, p)	Split x into set1 and set2 by keeping p% samples in set2
tokenization()	tokenize and input formatting function
select_tokenizer(x)	select x tokenizer to better train the data of Tweets
init_Tokenizer(x)	initialize the tokenizer with x pretrained model
addNomalization_L()	add normalization and make the text as Lowercase
encode_plus()	Encode_plus methods
indexes add(x,y)	add x and y tokens and return indexes by referencing tokenizer vocabulary

truncate(y)	truncate sentences to y length
create_attentionMask()	create attention mask
return_input_id()	returns token ids representing the tokenised sequence
return_attentionMask()	returns which tokens to attend to and which ones to ignore

Following algorithm includes the steps for Pre-processing:

1	Load(x)
2	Clean $(x, y_1, y_2 y_i)$
3	Split(x, set1, set2, p)
4	tokenization()
5	select_tokenizer(x)
6	init_Tokenizer(x)
7	addNomalization_L()
8	encode_plus()
	{
9	i = split_text();
10	j = add(CLS, SEP);
11	truncate(max_length)
12	create_attentionMask();
13	k= (input_ids[], attention_masks[])
14	return k;
	}
15	indexes add(x,y)
16	truncate(y)
17	create_attentionMask()
18	return_input_id()
19	return_attentionMask()

After loading the dataset into the system, cleaning the dataset is an important step in preparing data for our model - *TweetTruth*. The dataset should be pre-processed to remove any irrelevant information. A pre-processing pipeline involves removing URLs, mentions, hashtags, emojis, and punctuations from the tweets. To normalize the data there is need to convert all text to lowercase and replace user mentions with a special token.

Since it is a classification problem there is need to remove the stop words (common words like "the", "and", "a", etc.) and perform *stemming or lemmatization* (reducing words to their base form) on the text data. These steps can help reduce the dimensionality of the data and improve the performance of *TweetTruth*. After cleaning the data, the *train_test_split* function in the scikit-learn library is used for splitting a dataset into two subsets: a training set and a testing set. This is followed by tokenization.

Tokenization is a crucial step in natural language processing tasks that involves converting text into a sequence of tokens for further analysis. The tokenization process is encapsulated in a function called tokenization, which takes two arguments, X_train and X_val. The code snippet begins selecting the BERT tokenizer through the by select_Tokenizer(BERT) function call. The chosen tokenizer is initialized with the "vinai/tweettruth base" pretrained model using the *init_Tokenizer* function. Subsequently, text normalization techniques are applied using the addNormalization_L() function. The select_tokenizer is specifically designed for tokenizing text from social media platforms. The tokenizer is then initialized with the pretrained TweetTruth model by specifying its name or model identifier, "vinai/tweettruth_base."

The TweetTruth algorithm uses a pre-processing technique called "text normalization" to clean and standardize the tweets before feeding them into the model. The text normalization process involves several steps such as removing URLs, mentions, hashtags, emojis, punctuations, and stop words. The text is then tokenized into words and converted into lowercase. The algorithm also uses subword tokenization to handle out-of-vocabulary (OOV) words. In subword tokenization, the words are broken down into smaller subwords based on their frequency of occurrence in the training data. The subwords are then used as input to the BERT model. This is followed by the encode_plus() function. The encode_plus() function is used in tokenization libraries to encode text sequences and generate the necessary inputs for training or inference with transformerbased models like BERT.

Tokenization and Input formatting functions:

1. To better train the data of tweets, we use a twitter specific Bert tokenizer - BERTweet.

2. Initialization of the tokenizer using the method AutoTokenizer.from_pretrained("vinai/tweettruth-

base"). Addition of normalization and make the text as lower case.

3. This will return the input ids and attention masks and to calculate the max length of the tweets, the encode method is used for the same and set the max length variable to the maximum length of the sentences after adding special tokens.

4. The above-mentioned Tokenization and input

formatting function is used to tokenize (X_train) and validation (X_val) data of tweets.

The encode_plus method is a 5-step process. The Figure 2 shows the steps which are used.



Fig 2. Steps performed by encode_plus()

The explanation of each of the steps of the encode_plus() function is given below:

i. tokens split_text()

The input text is split into individual tokens or sub-words. The splitting process depends on the specific tokenizer being used and the tokenization rules it follows. The tokenizer may split text based on whitespace, punctuation, or language-specific rules. The result is a sequence of tokens.

ii. indexes add(x,y)

Special tokens like [CLS] (start of the sequence) and [SEP] (separator between sentences or segments) are added to the tokenized sequence. These tokens are necessary for models like TweetTruth, as they provide important positional information.

iii. truncate(y)

If the tokenized sequence exceeds a maximum length specified by the model or task requirements, it may need to be truncated. Truncation involves removing tokens from the sequence to fit within the length limit.

iv. create_attentionMask()

An attention mask is a binary tensor that indicates which tokens should be attended to and which ones should be ignored during model training. The attention mask has the same length as the tokenized sequence and consists of 1s and 0s, where 1 indicates a valid token and 0 indicates a padded or ignored token.

v. create_input_id()

Input IDs are the tokenized sequence converted into a sequence of token IDs that correspond to the vocabulary of

the pre-trained model. Each unique token is mapped to a specific ID. The input IDs represent the input data that the model will process during training or inference.

vi. return_attentionMask()

Returning the Outputs: The encode_plus() function typically returns the following outputs:

a. input_ids: The token IDs representing the tokenized sequence.

b. attention_mask: The attention mask indicating which tokens to attend to and which ones to ignore.

Additional outputs may include token type IDs or segment IDs, depending on the model architecture and specific use case. These encoded inputs (input IDs, attention masks, etc.) can be passed to a transformer-based model for training, fine-tuning, or inference. The flow diagram for the Preprocessing step of TweetTruth is given in Figure 3.



Fig 3. Flow diagram for Pre-Processing Step of TweetTruth Model

3.2 Processing the Tokenized data using PyTorch DataLoader

PyTorch DataLoader is a utility in the PyTorch library that provides an efficient and convenient way to load and iterate over datasets during model training or evaluation. It allows us to handle various aspects of data loading, such as batching, shuffling, parallelizing data loading, and more. Table 2 lists the functions used in the algorithm and a short description for each.

Table 2. Functions for algorithm for processing the
tokenized data

Function	Description
create_maxlength()	Create max length from actual tweets
max(encode_plus())	use the encode method and set the maximum length variable to the maximum length of the sentences after adding special tokens
tokenization()	Assigning tokens to each word
isNotComplete(x)	return True till x is not completely trained/ validated). The function will keep on iterating until the dataset is completely trained / validated.
dataLoader()	Use pytorch dataloader class to save memory and boost the training speed
dataLoader conver()	Convert the tokenized train input_id , attention_masks, train labels to Dataloader class using RandomSampler to randomize drawing samples from dataset using batch size of 32 and return DataLoader object

Following algorithm includes the steps for processing the tokenized data:

Al	gorithm 2: Preprocessing of Tokenized
	data
1	create_maxlength()
2	max(encode_plus())
3	tokenization()
4	while \leftarrow (isNotComplete(x)) Do
	{
5	dataLoader()
6	dataLoader conver()
	}
7	create_maxlength()

- 8 tokenization()
- 9 isNotComplete(x)
- 10 dataLoader()

To efficiently process the tokenized data during training, we can convert it to a PyTorch DataLoader object. The DataLoader wraps the tokenized inputs (input IDs and attention masks) and labels, allowing us to iterate over minibatches of data during training.

RandomSampler is a sampling strategy provided by PyTorch. It shuffles the dataset to ensure that the samples drawn during training are randomized. This randomization helps us to reduce bias and correlation in the training process and promotes better model generalization.

The flow Diagram for Processing the Tokenized data using DataLoader TweetTruth is given in Fig 4.



Fig 4. Flow Diagram for Processing the Tokenized data using DataLoader TweetTruth Model

3.3 Creating the BERT Based Model Using Freezing BERT Layers and Employing a Feed-Forward Neural Network.

 Table 3. Functions for algorithm for creating the BERT

 based model

Function	Description
For Alg	gorithm 3
create_BERT (classifier class name, Hidden Size of BERT, hidden Size of Classifier, no of output classes)	Creates the BERT model by initializing the given parameters

init_model (x)	initialize the model by
_ 、,	loading the pretrained
	model x
init_ff(class name, no of	initializes a feed-forward
input layers, activation	(FF) layer, which is
lunctions, no of output	TweetTruth model for
layer)	classification tasks The
	'Sequential' function is
	used to define a sequential
	stack of layers, where 'I'
	represents the input size and
	activation function.
float forward(x y) \rightarrow take	extracts the last hidden state
x and y as parameters and	of the [CLS] token and
returns the probabilities.	provide it to the
	classification model and
	returns the probabilities.
forward(x,y)	take x and y as parameters
	and returns the probabilities
init_paramters	initialize the model
(opunitzer_name, learning_rate, default)	training steps as length of
lourining_rate, derault)	train_dataloader as no. of
	times as epochs. Define the
	learning rate schedular
	using transformers
	warmup()
get_linear_scedule_with	Define the learning rate
_warmup()	schedular using
	transformers
	warmup ()
classImbalance(x)	Use sklearn's compute
erussime unice(n)	class weight() using
	class_weight parameter to
	handle class imbalance
convert_cw_to_to()	Convert class weight to
	tensor objects to move them
	then to GPU
use_cw()	uses the computed class
	weights in model.

For Algorithm 4

model.train()	put the model into the training model
load_batch(GPU)	load batch to the GPU
perform _i (bert)	perform a forward pass of the Bert model
perform _{i-1} (bert)	perform a backward pass of the Bert model

computeLV()	convert the loss value
dump(x)	print out loss values after x no of batch
total(x[])	find total of the given set x
save_model()	save the current model
For Alg	gorithm 5
load_batch(GPU)	loads a batch of data onto the GPU for efficient computation.
$gradient_{i-1} = 0$	initializes the gradient to zero at the beginning of each step. This is typically done to clear any previously accumulated gradients.
perform(bert)	performs the forward pass and backward pass (backpropagation) through the model to compute the gradients. The variable LV[i] is assigned the computed loss value (LV) for the current step.
perform _{i-1} (bert)	performs a forward pass on the TweetTruth model with the updated parameters.
$P_{i+1} = P_i + \alpha$	updates the parameter values (P) based on the learning rate (α) and LR _{i+1} = L _i + β updates the learning rate (LR) based on the weight decay value (β).
$A_{LV} = total(LV[i])$	computes the total loss A_{LV} by summing the loss values for each step within the epoch
If ($VL_i < VL_{i-1}$)	statement checks if the loss value for the current epoch VL_i is less than the loss value for the previous epoch VL_{i-1} .
save_model()	If the condition is met (i.e., the loss decreases), save the model's parameters and state to a file.

Following algorithm includes the steps for creating the BERT based model:

Algorithm 3: Bert-based Model

1 create_BERT (classifier class name, Hidden Size of BERT, hidden Size of Classifier, no. of output classes)

- 2 init_model (x)
- 3 init_ff(class name, no of input layers, activation functions, no of output layer)
- 4 float forward(x,y)
- 5 init_paramters (optimizer_name, learning_rate, default)
- 6 classImbalance(x)
- 7 convert_cw_to_to()
- 8 use_cw()
- 9 create_BERT (classifier class name, Hidden Size of BERT, hidden Size of Classifier, no of output classes)
- 10 init_model (x)
- 11 init_ff(class name, no of input layers, activation functions, no of output layer)
- 12 forward(x,y) \rightarrow take x and y as parameters and returns the probabilities.
- 13 init_paramters (optimizer_name, learning_rate, default)
- 14 get_linear_scedule_with _warmup()
- 15 classImbalance(x)
- 16 convert_cw_to_to()
- 17 use_cw()

Algorithm 4

- 1 for \leftarrow i to n epoch do
- 2 model.train()
- 3 load_batch(GPU)
- 4 perform_i(bert)
- 5 perform $_{i-1}$ (bert)
- 6 computeLV()
- 7 dump(x)
- 8 total(x[])
- 9 save_model()

Inside the epoch loop, the following steps are performed:

a. model.train() sets the model in training mode, enabling the gradients to be computed and parameters to be updated during backpropagation.

b. for_each_step_and_batch() is a nested loop that iterates over each step and batch within an epoch.

Inside this loop, as also depicted in Fig 5, the following steps are performed

	Algorithm 5	
1	load_batch(GPU)	
2	$gradient_{i-1} = 0$	
3	perform(bert)	
4	$perform_{i-1}(bert)$	
5	$P_{i+1} = P_i + \alpha$	
6	$A_{LV} = total(LV[i])$	
7	$\mathrm{lf}(\mathrm{VL}_{i}<\mathrm{VL}_{i-1})$	
8	save_model().	

The flow diagram for Fine tuning the TweetTruth is given in Fig 5.



Fig 5. Flow Diagram for Fine tuning the TweetTruth Model

3.4 Evaluation of the Trained Model to Classify Different Types of Fake News

Table 4 lists the functions used in the algorithm for evaluating the BERT based model and a short description for each, while the entire flow of this step is given in Figure 6.

Table 4. Functions for algorithm for evaluating the BERT based model

Function	Description	
For Algorithm 6		
model.eval()	put model into the evaluation	
compute(AA)	compute average performance metric	
compute(VS)	computes a validation score	
train(BERT)	train the BERT based model	
predict()	Predict method defining	
save_model()	saving the trained model	
deploy(model,A)	deploy the model to do inference is A value is satisfactory	
For Algorithm 7		
load_batch(GPU)	loads a batch of data onto the GPU for efficient computation.	
compute(L)	computes the loss (L) for the batch. The loss function depends on the specific task being performed, such as cross-entropy loss for classification tasks.	
compute(l) computes the	average loss (l) for the entire evaluation set. It accumulates the losses computed in each batch.	
get(Pr)	gets the predicted probabilities (Pr) from the model for the batch inputs.	
compute(A)	computes a performance metric (e.g., accuracy, F1 score) using the predicted probabilities and the ground truth labels for the batch.	

Following algorithm includes the steps for evaluating the BERT based model:

Algorithm 6: Evaluation of Bert-based Model

- 1 model.eval()
- 2 compute(AA)
- 3 compute(VS)
- 4 train(BERT)
- 5 predict()
- 6 save_model()

7 deploy(model,A)

model.eval() sets the model in evaluation mode, where the model parameters are fixed and no gradients are computed. This mode is typically used during inference or evaluation to disable any dropout or batch normalization layers that are active during training.

Inside the batch loop of dataloader, the following steps are performed:

Algorithm 7			
1	load_batch(GPU)		
2	compute(L)		
3	compute(l)		
4	get(Pr		
5	compute(A)		
nside	the train(BERT) function, the following steps a		

Inside the train(BERT) function, the following steps are performed:

i. ne = 5 sets the number of epochs (ne) for training.

ii.*init_model(BERT)* initializes the BERT model for training.

iii.*train_start()* starts of the training process, which involves iterations over epochs and batches, similar to what we previously described in the training loop.

Inside the predict() function, the following steps are performed:

i.*compute_f(x)* computes x using forward pass.

ii.*softmax()* function is applied to the model's output probabilities (prob). Softmax is commonly used for multi-class classification problems to convert raw scores into probabilities.

iii.*model.eval()* puts the model into the evaluation mode.

iv.concatenate (x) concatenates the individual loss values (L[i]) into a single tensor or array.

v.generate(x,y()) generate x curve using y method

vi. *generate*(*AUC*,*y*()) Area under curve using y method

vii.generate(ROC,y()) Receiver Operating characteristics
curve using y method Inside the save_model() function,
the following steps are performed:

i. *generate(A, Predict())* generates evaluation metric A by calling the Predict() function.

ii. *generate(AUC, Predict())* generates another evaluation metric, the Area Under the Curve (AUC), by calling the Predict() function.

iii. *generate(ROC, Predict())* generates the Receiver Operating Characteristic (ROC) curve by calling the Predict() function.

iv. *deploy(model, A)* deploys the trained TweetTruth model to do inference if A value is satisfactory.



Fig 6. Evaluating the TweetTruth Model

4. Result

In the pursuit of mitigating the proliferation of deceptive news on the social media platform X (formerly known as Twitter), we implemented an innovative strategy centered around BERT (Bidirectional Encoder Representations from Transformers) technology. Our approach involved meticulous data acquisition, preprocessing, fine-tuning, and optimization of the BERT model to address the dissemination of misleading or fabricated content. The "socialmedia-disaster-tweets-DFE" dataset (27) is used to test our model performance.

The "socialmedia-disaster-tweets-DFE" dataset is a valuable resource comprising a collection of tweets from various social media platforms, particularly focusing on content related to disasters and emergencies. This dataset has been curated to include tweets that contain information, discussions, or reports related to disasters such as natural calamities, accidents, or crises of different magnitudes. It encompasses a wide array of textual data, encompassing user-generated content that can include firsthand accounts, reactions, and discussions pertaining to ongoing or recent disasters. Table 6 shows the performance metrics of the proposed TweetTruth model while Figure 7 represents the Receiver Operating Characteristic Curve for the model.

Table 6. Performance metrics of TweetTruth

Metric	Value
Accuracy	0.8441
Precision	0.693
Recall	0.77
F1 Score	0.73



Fig 7. ROC of the model

Fig 7 shows a comparison of the various metrics between the proposed TweetTruth model and existing state of the art models for the socialmedia-disaster-tweets-DFE. TweetTruth surpasses them all with an accuracy of 84.4 %, precision of 69.3%, Recall score of 77% ad an F1 score of 73%.



Fig 8. Comparison of TweetTruth with state-of-the-art models

These results underscore the efficacy of our BERT-based approach in combating the spread of deceptive news on social media, particularly on platform X. While the model's accuracy and AUC are promising, the F1 score reflects a balanced trade-off between precision and recall. Our research contributes to the development of robust natural language processing techniques aimed at fostering a more dependable information environment in the digital era. These findings hold significant implications for enhancing the integrity and credibility of news dissemination on social media platforms.

5. Conclusion and Future Work

The paper proposes a robust approach for detecting fake news on Twitter using a BERT-based deep learning algorithm. By leveraging the power of BERT's contextual understanding and linguistic features, the algorithm effectively distinguished between genuine and fake news tweets. The study included a comprehensive literature survey, highlighting the effectiveness of BERT-based models and other machine learning techniques in addressing the challenge of fake news detection on social media platforms. The proposed algorithm demonstrated high accuracy rates and outperformed previous approaches in identifying fake news.

The implications of this research are significant, as accurate fake news detection on Twitter contributes to preserving information integrity and promoting informed decisionmaking. By mitigating the spread of misinformation, the algorithm helps protect individuals from potential harm caused by false information. Future research directions may include refining the algorithm by incorporating additional contextual features, exploring multi-modal approaches, and addressing the challenges of evolving fake news techniques. Overall, this research contributes to the development of robust mechanisms for combating fake news on Twitter, leading to a more trustworthy and reliable social media environment.

References

- Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. Science, 359(6380), 1146-1151.
- [2] Shao, C., Ciampaglia, G. L., Varol, O., Yang, K. C., Flammini, A., & Menczer, F. (2018). The spread of low-credibility content by social bots. Nature communications, 9(1), 1-10.
- [3] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Bidirectional Encoder Representations from Transformers. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT), 4171-4186.
- [4] H. C. Jwa, H. Kim, and J. Park, "FakeBERT: Fake news detection in social media with a BERT-based deep learning approach," Information Sciences, vol. 521, pp. 118-134, 2020.
- [5] N. Bounaama and M. Abderrahim, "Classifying COVID-19 related tweets for fake news detection

and sentiment analysis with BERT-based models," Applied Sciences, vol. 13, no. 12, p. 6275, 2023.

- [6] S. Singh, P. Shukla, and H. Kaur, "Fake or real news about COVID-19? Pretrained transformer model to detect potential misleading news," Engineering Applications of Artificial Intelligence, vol. 116, p. 104113, 2022.
- [7] J. Zhang, Z. Wang, and Y. Liu, "BERT-Tweet: A pre-trained language model for English Tweets," arXiv preprint arXiv:2001.08308, 2020.
- [8] R. Kumar, A. Verma, and N. Singh, "Fake news detection on Twitter using natural language processing techniques," Information Processing & Management, vol. 57, no. 2, pp. 367-384, 2020.
- [9] A. Verma, N. Singh, and R. Kumar, "A deep learning approach to fake news detection on Twitter," Expert Systems with Applications, vol. 114, pp. 137-147, 2018.
- M. Imran, M. Mohsin, and J. Qadir, "Detecting fake news on Twitter using machine learning techniques," Information Processing & Management, vol. 53, no. 4, pp. 779-791, 2017.
- [11] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "Fake news detection on Twitter with feature engineering and machine learning," arXiv preprint arXiv:1609.08253, 2016.
- [12] F. Wang, J. Xu, Z. Li, and H. Zhang, "Toward Deeper Understanding of Neural Networks: The Power of Initialization and a Dual View on Expressivity," in Proceedings of the 35th International Conference on Machine Learning (ICML), 2018.
- [13] Z. Ma and A. Sun, "Detecting Rumors from Microblogs with Recurrent Neural Networks," in Proceedings of the 25th International Conference on World Wide Web (WWW), 2016.
- [14] S. Wu, D. Yang, and N. Xu, "Who Falls for Fake News? The Roles of Analytic Thinking, Motivated Reasoning, Political Ideology, and Bullshit Receptivity," in Journal of Applied Research in Memory and Cognition, 2020.
- [15] C. Castillo, M. Mendoza, and B. Poblete, "Information Credibility on Twitter," in Proceedings of the 20th International Conference on World Wide Web (WWW), 2011.
- [16] N. Ruchansky, S. Seo, and Y. Liu, "CSI: A Hybrid Deep Model for Fake News Detection," in Proceedings of the 2017 ACM on Conference on Information and Knowledge Management (CIKM), 2017.

- [17] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "Fake news detection on Twitter with feature engineering and machine learning," arXiv preprint arXiv:1609.08253, 2016.
- [18] H. Ahmed, I. Traore, and S. Saad, "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, pp. 127–138, 2017.
- [19] A. Reema, A. K. Kar, and P. Vigneswara Ilavarasan, "Detection of spammers in twitter marketing: a hybrid approach using social media analytics and bio-inspired computing," in Information Systems Frontiers, vol. 20, no. 3, pp. 515–530, 2018.
- [20] K. Shu, S. Wang, and H. Liu, "Beyond news contents: The role of social context for fake news detection," in Proceedings of the twelfth ACM international conference on web search and data mining, pp. 312–320, ACM, 2019.
- [21] D. Munandar, A. Arisal, D. Riswantini, and A. F. Rozie, "Text classification for sentiment prediction of social media dataset using multichannel convolution neural network," in 2018 International conference on computer, control, informatics and its applications (IC3INA), IEEE, pp. 104–109, 2018.
- [22] A. Gupta and P. Kumaraguru, "Credibility ranking of tweets during high impact events," in Proceedings of the 1st Workshop on Privacy and Security in Online Social Media, ser. PSOSM '12, New York, NY, USA, ACM, pp. 2:2–2:8, 2012.
- [23] S. Mohd Shariff, X. Zhang, and M. Sanderson, "User perception of information credibility of news on Twitter," in Proceedings of the 36th European Conference on IR Research on Advances in Information Retrieval - Volume 8416, ser. ECIR 2014, New York, NY, USA, Springer-Verlag New York, Inc., pp. 513–518, 2014.
- [24] S. Sikdar, S. Adali, M. Amin, T. Abdelzaher, K. Chan, J. H. Cho, B. Kang, and J. O'Donovan, "Finding true and credible information on Twitter," in 17th International Conference on Information Fusion (FUSION), July 2014, pp. 1–8.
- [25] Author(s), "Title of the Webpage," Zenodo,[Online]. Available: https://zenodo.org/. [Accessed: August 16, 2023].
- [26] W. McKinney, "Pandas: a foundational Python library for data analysis and statistics," IEEE Transactions on Visualization and Computer Graphics, vol. 26, no. 1, pp. 55-66, 2020. DOI: 10.1109/TVCG.2019.2930765.

International Journal of Intelligent Systems and Applications in Engineering

- [27] "Disasters on Social Media," Kaggle, [Online]. Available: https://www.kaggle.com/datasets/jannesklaas/disast ers-on-social-media. [Accessed: August 16, 2023].
- [28] C. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer, "Online Passive-Aggressive Algorithms," Journal of Machine Learning Research, vol. 7, pp. 551–585, 2006.
- [29] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, "Language models are unsupervised multitask learners," OpenAI Blog, vol. 1, no. 8, p. 9, 2019.
- [30] Archana Nanade, Dr. Amit Jain, Dr. Prateek Srivastava, Shweta Lalwani Hod, "Machine Learning Based Fake News Detection Using Natural Language Processing", IJAST, vol. 29, no. 08, pp. 5988 - 6003, Nov. 2.
- [31] Martinez, M., Davies, C., Garcia, J., Castro, J., & Martinez, J. Machine Learning-Enabled Quality Control in Engineering Manufacturing. Kuwait Journal of Machine Learning, 1(2). Retrieved from http://kuwaitjournals.com/index.php/kjml/article/vie w/122
- [32] Thota, D. S. ., Sangeetha, D. M., & Raj, R. . (2022). Breast Cancer Detection by Feature Extraction and Classification Using Deep Learning Architectures. Research Journal of Computer Systems and Engineering, 3(1), 90–94. Retrieved from https://technicaljournals.org/RJCSE/index.php/journ al/article/view/48