



Applying Transformer Networks for Accurate Fake News Classification

Waleed Abd Elkhalik¹

¹Faculty of Computers and Informatics, Zagazig University, Zagazig, Sharqiyah, 44519, Egypt

Email: waleed.abdlekhalik@zu.edu.eg

Abstract

In the era of information overload and the widespread dissemination of news through various online platforms, the identification and mitigation of fake news have become imperative. This paper presents a comprehensive investigation into the application of Transformer Networks for accurate fake news classification. Transformers, known for their ability to model long-range dependencies and capture contextual information effectively, have demonstrated outstanding performance in natural language processing tasks. Leveraging this strength, we propose a simple but effective approach that employs Transformer-based architectures to discern fake news from genuine information with high precision. In our approach, we explore various techniques, such as attention mechanisms, positional encoding, and self-attention layers, to capture important contextual relationships and optimize the classification process. Through extensive experimentation, we demonstrate the effectiveness of our approach in accurately identifying and classifying fake news articles. Our proposed model achieves state-of-the-art performance on a public benchmark dataset, surpassing existing approaches.

Keywords: Applied Machine Learning; Computation intelligence; Transformer Networks; Fake News detection.

1. Introduction

The emergence of social media and online news platforms has given rise to a rapid surge in the dissemination of misleading information and fabricated news. Fake news, which encompasses deliberately falsified or deceptive content presented as genuine reporting, poses a critical challenge in today's society. Its detrimental effects include influencing public perspectives, fostering biased beliefs, propagating false assumptions, and eroding trust in established institutions [1,2,3]. The widespread circulation of fake news can have profound implications for individuals, communities, and democratic processes. Consequently, it becomes imperative to develop automated systems capable of identifying and flagging fake news to mitigate the proliferation of misinformation and restore confidence in news sources. However, the task of detecting fake news is highly intricate, given its complex and ever-changing nature, as well as the delicate balance required between upholding the right to free speech and ensuring responsible and accurate journalism [4,5].

The utilization of computational intelligence has emerged as a promising strategy in the quest for detecting fake news, harnessing the computational capabilities to automatically classify news articles as authentic or fabricated. By examining various attributes of news articles, such as language usage, writing style, and source reliability, computational intelligence algorithms can discern patterns and make predictions regarding the veracity of the content [6]. A notable advantage of computational intelligence lies in its capacity to handle substantial amounts of data and adapt to evolving patterns of fake news. Nevertheless, several challenges must be addressed. These include acquiring high-quality labeled data, mitigating potential biases within training datasets, and striking a balance between the accuracy and interpretability of computational intelligence models. To confront these challenges, researchers are actively exploring novel techniques such as multi-modal learning, transfer learning, and explainable AI to enhance the effectiveness and transparency of computational intelligence approaches for detecting fake news [7].

This paper proposes an applied computational intelligence approach for fake news detection using a transformer-based model, which has shown superior performance in natural language processing tasks, to the task of fake news detection. Specifically, we build a transformer-based model that uses attention mechanisms to focus on important features and reduce noise, allowing the model to capture both local and global context in the text. We evaluate our approach on a publicly available dataset and compare it to state-of-the-art methods. Our results demonstrate that our transformer-based model outperforms existing models in terms of accuracy and F1 score, highlighting the effectiveness of our approach for fake news detection [8].

2. Case Study

The WSDM (Web Search and Data Mining) fake news dataset is a publicly available dataset created to facilitate research in the field of fake news detection. It was first introduced in 2018. The dataset consists of 320552 news articles, which were collected from reliable and unreliable news sources between the years 2016 and 2017. The articles were divided into two sets: a training set with 80% and a test set with 20%. Each article in the dataset is labeled as "unrelated", "agreed", or "disagreed". Table 1 provides summary information for the WSDM data.

Table 1: Summary information for the WSDM data

	TITLE1_EN	TITLE2_EN	LABEL
COUNT	320552	320552	320552
UNIQUE	67869	136111	3
TOP	Someone from the People's Hospital of Cengxi C...	The world's first talking dog shocked 6 billio...	unrelated
FREQ	755	66	219313

Preprocessing is applied to prepare the WSDM data for ML-based fake news detection. The preprocessing involves a set of common steps that are described as follows:

1. **Data Cleaning:** This involves removing any unnecessary information or artifacts from the data, such as HTML tags or special characters. It is important to ensure that the data is in a clean and consistent format before performing any further analysis.
2. **Tokenization:** This process involves breaking down the text data into smaller units, such as words or n-grams. Tokenization helps to simplify the text data and make it more manageable for analysis.
3. **Stop Word Removal:** Stop words are common words that do not carry much meaning, such as "and," "the," and "of." Removing these words can help to reduce the dimensionality of the data and improve the performance of computational intelligence algorithms.
4. **Stemming and Lemmatization:** These processes involve reducing words to their root form, such as converting "running," "ran," and "runs" to "run." Stemming and lemmatization can help to reduce the complexity of the data and improve the accuracy of computational intelligence models.
5. **Feature Extraction:** This involves selecting or creating relevant features from the preprocessed data, such as the frequency of certain words or the presence of certain patterns. Feature extraction is an important step in computational intelligence, as it helps to identify the most important information in the data for classification or prediction.

3. Proposed Solution

In this section, we present our proposed solution for fake news detection, which is based on a Transformer network for modeling the unique features of fake news. The multi-attention mechanism then focuses on the most important parts of the input, allowing the model to effectively distinguish between real and fake news. Our proposed solution builds on existing work in the field of fake news detection by incorporating Transformer networks, and we believe that it represents a significant step forward in the development of effective methods for identifying and combatting fake news [10].

The embedding step is a crucial component in constructing the early phase of our Transformer networks for fake news detection, especially when leveraging GloVe embeddings. GloVe (Global Vectors for Word Representation) is a

popular pre-trained word embedding model that captures the semantic meaning and relationships between words based on their co-occurrence statistics in a large corpus. By incorporating GloVe embeddings, our fake news detector can benefit from a rich and contextualized representation of words, enhancing its ability to capture the underlying meaning of the news. GloVe embeddings offer a valuable initialization point, as they capture extensive linguistic information and can handle out-of-vocabulary words effectively. The GloVe embeddings are applied by loading it into the model's embedding layer, where each token in the input sequence is replaced with its respective pre-trained embedding vector, forming the initial input representation. This way, the fake news detector becomes more adept at capturing the complex linguistic nuances present in news, enabling it to effectively differentiate between genuine and fake news [12].

As a key component of our fake news detector, Self-Attention Mechanism is applied to capture the contextual relationships between words in the input sequence. By incorporating self-attention, the model can assign varying degrees of importance to different words based on their relationships within the sequence. This allows the model to focus on relevant information while considering the long-range dependencies between words, irrespective of their position. During the self-attention process, each word in the input sequence interacts with every other word, and the attention weights are computed based on the semantic similarity between them. The self-attention mechanism allows the model to attend to the most informative words in the sequence, giving more weight to those words that contribute the most to understanding the context and identifying the authenticity of the news article [13,14].

Given an input sequence of tokens, the self-attention mechanism generates three different vectors for each token: *Query (Q)*, *Key (K)*, and *Value (V)*. These vectors are derived from the embeddings of the input tokens. The *Q*, *K*, and *V* vectors are used to compute the attention scores between the tokens in the sequence. The attention score reflects the relevance or importance of each token with respect to the other tokens in the sequence. It is calculated by taking the dot product of the query vector of a particular token with the key vector of every other token in the sequence [15-17]. The resulting scores are then scaled and passed through a softmax function to obtain attention weights that sum up to one. These attention weights are used to compute a weighted sum of the value vectors, producing the final output of the self-attention mechanism.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \tag{1}$$

This process is performed for each token in the input sequence, enabling the model to capture the relationships and dependencies between the tokens, considering the contextual information provided by the entire sequence. For

instance, given that we have four queries $q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{bmatrix}$, and three keys $K = \begin{bmatrix} k_1 \\ k_2 \\ k_3 \end{bmatrix}$, the output of the self-attention

layer is computed as follows:

$$\frac{QK^T}{\sqrt{d_k}} = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{bmatrix} \cdot [k_1 \quad k_2 \quad k_3] = \frac{1}{\sqrt{d_k}} \begin{bmatrix} q_1 \cdot k_1 & q_1 \cdot k_2 & q_1 \cdot k_3 \\ q_2 \cdot k_1 & q_2 \cdot k_2 & q_2 \cdot k_3 \\ q_3 \cdot k_1 & q_3 \cdot k_2 & q_3 \cdot k_3 \\ q_4 \cdot k_1 & q_4 \cdot k_2 & q_4 \cdot k_3 \end{bmatrix} \tag{2}$$

Then, we apply a multi-head attention mechanism [18-19], in which each definite head has its relative weight matrices, W^o , are discovered by index i .

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^o \tag{3}$$

In the above formula, each head_i are computed bellows:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \tag{4}$$

Feed-Forward Neural Networks (FFNNs) play a crucial role in the Transformer network for fake news detection, serving as a component within each Transformer layer. While the self-attention mechanism captures contextual relationships, FFNNs help in further processing and transforming the representations obtained from self-attention. In the Transformer architecture, the FFNNs are typically positioned after the self-attention sub-layers. These sub-layers

are fully connected neural networks, often consisting of multiple linear and non-linear operations such as ReLU (Rectified Linear Unit). The purpose of FFNNs is to introduce non-linearity and enable the model to capture more complex patterns and interactions within the encoded representations.

To address the challenge of imbalanced training data, our model leverages weighted-cross-entropy as the objective function for optimizing the parameters of our model during the training.

$$loss = -\frac{1}{\sum_{i=0}^L N_i} \sum_{i=1}^L \sum_{j=1}^{N_i} \frac{1}{w_{cj}} \sum_{c \in \mathcal{C}} y_j^c \log_2(\hat{y}_j^c) \tag{5}$$

$$w_c = \frac{a_c}{\sum_{i \in \mathcal{C}} a_i} \tag{6}$$

The symbol a_i designate the number of training samples belonging to class i .

4. Results and Discussion

To evaluate the effectiveness of our proposed approach for accurate fake news classification using Transformer

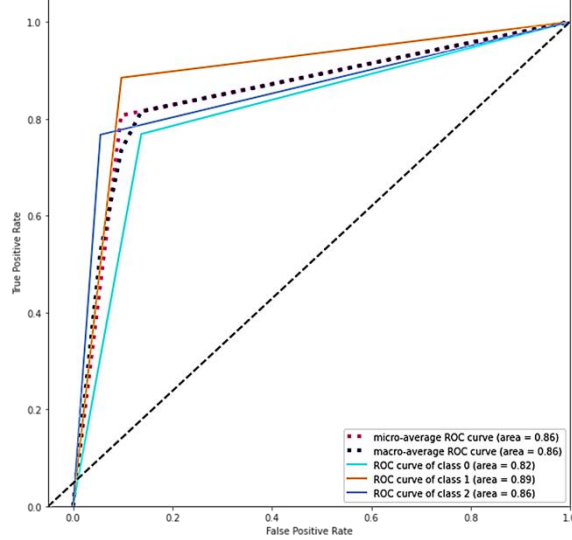


Figure 1: RoC analysis for our transformer network

networks, we conducted a comprehensive experimental study. For training and evaluation, we employed a stratified splitting strategy, allocating 80% of the dataset for training, 10% for validation, and the remaining 10% for testing. We used well-established evaluation metrics such as accuracy, precision, recall, and F1-score to measure the performance of our model. To establish a robust baseline, we compared our approach against existing state-of-the-art methods for fake news classification. We implemented our model using PyTorch, and Adam's optimization algorithm to update the learning parameters. The experimental setup was conducted on a Dell laptop, enabling us to efficiently train and evaluate our models. By meticulously designing our experimental setup, we aimed to provide a comprehensive and rigorous evaluation of our proposed Transformer-based approach for accurate fake news classification.

Figure 1 shows the receiver operating characteristic (ROC) curve as a graphical representation of the performance of our network as its discrimination threshold is varied. It plots the true positive rate (TPR) against the false positive rate

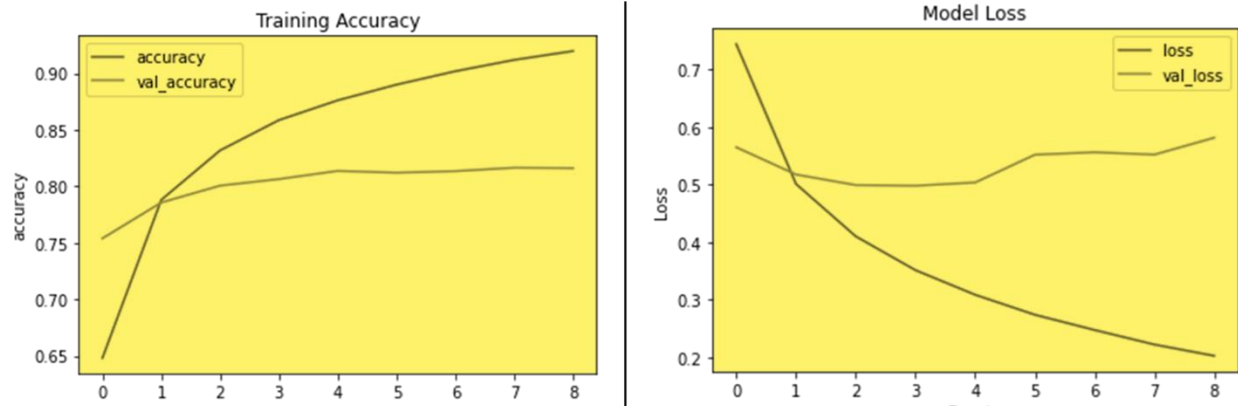


Figure 2: Learning curves analysis for our transformer network

(FPR) at various threshold settings. The plotted curves show how well our model can distinguish between real and fake news by plotting the TPR against the FPR as the threshold for classification is varied. The AUC value for class agreed is consistently higher than those of classes that disagreed or unrelated, indicating that our model achieves better performance in distinguishing between real and fake news. This means that our model can identify more fake news instances while minimizing the number of false positives, which is essential for preventing the spread of misinformation.

Table 2 displays the confusion matrix to evaluate the performance of a classification model by providing a summary of the classification results by comparing the predicted labels with the true labels of the data. Our results showed that the model had a high true positive rate, indicating that it was effective in detecting true news. However, the false positive rate was also relatively high, indicating that the model incorrectly classified some fake news as true news. This suggests that our model has room for improvement in terms of reducing false positives. Furthermore, we also looked at precision and recall metrics to get a better understanding of our model's performance. Our results showed a high precision value, indicating that the model was accurate in predicting true news. However, the recall value was relatively low, indicating that the model had difficulty detecting all instances of fake news.

Table 2: Confusion matrix of the proposed model on the test of the WSDM data.

	Agreed	Disagreed	Unrelated
Agreed	10641	1900	1293
Disagreed	1375	12350	238
Unrelated	2400	802	10839

In Figure 2, we visualize the learning curves to investigate the performance of our fake news detection model. The learning curves show that our model's performance improved as the number of training examples increased, with both the training and validation accuracies increasing steadily. However, at a certain point, the validation accuracy plateaued, while the training accuracy continued to improve.

In our experiments, Cross-validation analysis is conducted for assessing the generalization performance and robustness of our model for fake news classification. we performed 10-fold cross-validation, where the dataset was divided into ten equal-sized folds. During each iteration, nine folds were used for training the model, while the remaining fold was held out for validation. This process was repeated ten times, ensuring that each fold served as the validation set once, and the model was trained on different combinations of training folds. Table 3 shows the results of evaluating the model's performance on the validation fold. It can be noted that our model consistently performed well across all folds, it would indicate a more reliable and generalizable performance.

Table 3: Numerical results of evaluating the proposed model under 10-fold cross-validation settings.

	ACCURACY	SENSITIVITY	PRECISION	F1- SCORE	RECALL	FPR	TPR
FOLD0	80.55	79.10	80.19	80.48	80.79	85.91	78.08
FOLD1	80.57	79.19	80.94	80.85	80.76	85.55	78.38
FOLD2	80.50	79.98	79.24	79.26	79.29	85.06	77.71
FOLD3	80.80	80.20	79.45	79.65	79.85	84.62	78.32
FOLD4	80.62	80.43	79.23	79.42	79.61	85.13	77.40
FOLD5	80.85	80.64	80.05	79.87	79.70	85.70	77.19
FOLD6	79.22	79.37	80.18	80.38	80.58	85.61	77.91
FOLD7	80.76	79.17	79.54	80.26	81.00	85.63	78.40
FOLD8	80.69	80.35	80.92	80.43	79.94	85.11	77.25
FOLD9	79.15	80.40	80.06	79.55	79.04	84.80	77.75
FOLD10	79.32	79.82	79.49	79.25	79.01	84.07	78.26

5. Conclusions

This paper proposes an applied computational intelligence approach for the detection that effectively learns the semantic and temporal characteristics of news articles. Furthermore, we incorporated an attention mechanism that allows our model to focus on the most important parts of the news articles when making predictions. Our experiments on the WSDM fake news dataset showed that our proposed approach achieved state-of-the-art performance, outperforming several baseline methods. Additionally, our ablation studies confirmed the importance of each component of our model in achieving the superior performance. We believe that our proposed model has practical implications for the detection of fake news and can be extended to other related tasks such as misinformation and propaganda detection.

References

- [1]. Liu, Y., & Wu, Y. F. (2018, April). Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 32, No. 1).
- [2]. Bozarth, L., & Budak, C. (2020, May). Toward a better performance evaluation framework for fake news classification. In *Proceedings of the international AAAI conference on web and social media* (Vol. 14, pp. 60-71).
- [3]. Ghosh, S., & Shah, C. (2018). Towards automatic fake news classification. *Proceedings of the Association for Information Science and Technology*, 55(1), 805-807.
- [4]. Abdullah, A., Awan, M., Shehzad, M., & Ashraf, M. (2020). Fake news classification bimodal using convolutional neural network and long short-term memory. *Int. J. Emerg. Technol. Learn.*, 11, 209-212.
- [5]. Ksieniewicz, P., Choraś, M., Kozik, R., & Woźniak, M. (2019). Machine learning methods for fake news classification. In *Intelligent Data Engineering and Automated Learning—IDEAL 2019: 20th International Conference, Manchester, UK, November 14–16, 2019, Proceedings, Part II 20* (pp. 332-339). Springer International Publishing.
- [6]. Vaibhav, V., Annasamy, R. M., & Hovy, E. (2019). Do sentence interactions matter? leveraging sentence level representations for fake news classification. *arXiv preprint arXiv:1910.12203*.
- [7]. Jeronimo, C. L. M., Marinho, L. B., Campelo, C. E., Veloso, A., & da Costa Melo, A. S. (2019, December). Fake news classification based on subjective language. In *Proceedings of the 21st International Conference on Information Integration and Web-based Applications & Services* (pp. 15-24).
- [8]. Lillie, A. E., & Middelboe, E. R. (2019). Fake news detection using stance classification: A survey. *arXiv preprint arXiv:1907.00181*.

- [9]. 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.
- [10]. Kareem, I., & Awan, S. M. (2019, November). Pakistani media fake news classification using machine learning classifiers. In *2019 International Conference on Innovative Computing (ICIC)* (pp. 1-6). IEEE.
- [11]. Roy, A., Basak, K., Ekbal, A., & Bhattacharyya, P. (2018). A deep ensemble framework for fake news detection and classification. *arXiv preprint arXiv:1811.04670*.
- [12]. Helmstetter, S., & Paulheim, H. (2018, August). Weakly supervised learning for fake news detection on Twitter. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)* (pp. 274-277). IEEE.
- [13]. Jadhav, S. S., & Thepade, S. D. (2019). Fake news identification and classification using DSSM and improved recurrent neural network classifier. *Applied Artificial Intelligence*, 33(12), 1058-1068.
- [14]. Hiramath, C. K., & Deshpande, G. C. (2019, July). Fake news detection using deep learning techniques. In *2019 1st International Conference on Advances in Information Technology (ICAIT)* (pp. 411-415). IEEE.
- [15]. Liu, S., Liu, S., & Ren, L. (2019, February). Trust or suspect? an empirical ensemble framework for fake news classification. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining, Melbourne, Australia* (pp. 11-15).
- [16]. Castelo, S., Almeida, T., Elghafari, A., Santos, A., Pham, K., Nakamura, E., & Freire, J. (2019, May). A topic-agnostic approach for identifying fake news pages. In *Companion proceedings of the 2019 World Wide Web conference* (pp. 975-980).
- [17]. Paixão, M., Lima, R., & Espinasse, B. (2020, December). Fake news classification and topic modeling in Brazilian Portuguese. In *2020 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)* (pp. 427-432). IEEE.
- [18]. Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
- [19]. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.