

An adaptive approach to detecting fake news based on generalized text features

Andrii Shupta¹, Olexander Barmak¹, Adam Wierzbicki² and Tetiana Skrypnyk¹

¹ Khmelnytskyi National University, Institutes 11 st. 29016, Khmelnytskyi, Ukraine

² Polish-Japanese Academy of Information Technology, Koszykowa 86 st. 02-008, Warsaw, Poland

Abstract

Fake news has become a serious problem in recent years as they can quickly spread through social media and other online platforms. Various methods and materials can be used to detect fake news. One approach involves analyzing the content of the news, including the text and accompanying images or videos. Another approach involves considering the social context in which the news is spread, such as the news source and the mood of the people sharing them. An adaptive approach for detecting fake news using Natural Language Processing is presented in this work. It is proposed to use a feature vector constructed from generalized characteristics of news texts. The possibility of expanding the feature vector and training data sets to adapt the classifier to new types and types of fake news is also proposed. The experimental results presented qualitatively (visual analytics) and quantitatively (statistical metrics) demonstrate the ability of the proposed approach to detect fake news with sufficient quality (90%). Overall, the research aims to contribute to the development of a reliable and accurate system for detecting fake news, which may have important consequences for addressing this problem in modern society.

Keywords

Fake news, Fake news detection, Natural Language Processing

1. Introduction

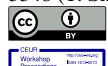
Fake news has become a serious problem in modern society, as it can quickly spread through social media and other online platforms, influencing people's thoughts and beliefs. Detecting fake news has become an important task that requires the use of various methods and techniques to accurately identify false or misleading information.

Social media is a primary means of news consumption, especially for younger individuals, but as the popularity of consuming news on social media platforms increases, so does the prevalence of misinformation, including false information and unsupported claims. Various methods based on text and social context have been developed to identify fake news on social media, but recent studies have explored the limitations and weaknesses of these fake news detectors. [1].

There are various social media platforms available to users, enabling them to post and share news online. These platforms lack verification measures for users and their posts, leading to the spread of false information by some users. Such misinformation can include propaganda targeted at individuals, society, organizations, or political parties. Due to the sheer volume of content, it is challenging for humans to detect all instances of fake news, highlighting the need for automated machine learning classifiers. [2].

Fake news detection methods are commonly trained on data that is available at the time of training, which may not be applicable to future events. This is because many of the labeled samples used for training on verified fake news may become outdated quickly as new events emerge. [3].

COLINS-2023: 7th International Conference on Computational Linguistics and Intelligent Systems, April 20–21, 2023, Kharkiv, Ukraine
EMAIL: andrii.shupta@gmail.com (A. Shupta); alexander.barmak@gmail.com (O. Barmak); adamw@pjwstk.edu.pl (A. Wierzbicki);
tkskripnik1970@gmail.com (T. Skrypnyk)
ORCID: 0009-0000-9771-5579 (A. Shupta); 0000-0003-0739-9678 (O. Barmak); 0000-0003-0075-7030 (A. Wierzbicki); 0000-0002-8531-5348 (T. Skrypnyk)



© 2023 Copyright for this paper by its authors.
Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).
CEUR Workshop Proceedings (CEUR-WS.org)

In the study, an adaptive approach to detecting fake news is proposed, based on a transparent, interpreted feature vector constructed from generalized characteristics of news texts. The adaptability of the approach lies in the ability to supplement the feature vector with new characteristics and to build a set of classifiers on different training sets.

The contributions of the article are as follows:

- an adaptive approach to fake news detection is proposed based on a feature vector constructed from generalized content characteristics;
- the ability of the proposed approach to detect fake news with acceptable values of statistical metrics is demonstrated.

The structure of the article is as follows: Section 2. Related works provides an overview and analysis of modern approaches to fake news detection and formulates the research goal. Section 3. Methods and Materials describes the proposed adaptive approach to fake news detection. Section 4. Results and discussion presents the research results, including visual and numerical values of statistical metrics, their correlation with similar research, and the confirmation of the proposed approach's ability to detect fake news. The further prospects of the proposed approach are discussed. Finally, the conclusions are presented.

2. Related work

Numerous works have been done to detect fake news using different techniques and methods. In [4], the authors proposed a novel method for detecting fake news by combining various features, including text and user-based features, and using deep learning models. The fundamental algorithms used in their study are an extension of traditional Convolutional Neural Networks (CNNs) to graphs. This enables the combination of dissimilar types of data such as content, user profile and activity, social graph, and news propagation. They achieved an accuracy of 92.7%.

Another study [5] focused on using linguistic features. Their study utilized a dataset comprised of two datasets containing an equal number of true and fake news articles related to politics. To extract linguistic and stylometric features, text fields from the dataset were utilized, and a bag of words TF and BOW TF-IDF vector were generated. A variety of machine learning models, including bagging and boosting methods, were then applied to achieve the highest level of accuracy.

In study [6], two machine learning algorithms were evaluated using word n-grams and character n-grams analysis for fake news detection. The experimental results showed that character n-grams combined with Term-Frequency-Inverted Document Frequency (TF-IDF) achieved better performance, with a Gradient Boosting Classifier achieving an accuracy of 96%.

Finally, in [7], the authors of this article proposed a theory-driven model to detect fake news, which examines news content at different levels, including the lexicon, syntax, semantics, and discourse. They used well-established theories in social and forensic psychology to represent news at each level and conducted fake news detection within a supervised machine learning framework. As an interdisciplinary study, their work aims to explore potential patterns in fake news, improve interpretability in fake news feature engineering, and investigate the relationships between fake news, deception/disinformation, and clickbaits.

Based on the analysis of related work, various weaknesses in the approaches can be identified. One of them is the inadequate quality of the data on which the model is based. If the model is trained on incorrect or insufficient data, it may classify news incorrectly.

Another factor is the speed at which news spreads on the Internet. Fake news can quickly gain popularity and spread faster than any model can detect them. It is also important to consider that fake news may contain some truthful information, making their detection more difficult.

Yet another reason is the changing technologies and approaches to creating fake news. As new technologies emerge over time that allow for more convincing fake news, models created to detect previous versions of fake news may be ineffective. It is also important to consider that most approaches to detecting fake news are based on machine learning, which can be vulnerable to attacks by malicious actors. For example, malicious actors can train the model to classify a certain type of news as fake by changing the content of the news.

Therefore, the aim of this work is to propose an approach that can adapt to the changing nature of fake news. The approach should retrain on new data, use previous results, and improve the accuracy of detecting fake news. Additionally, the approach should allow for expanding the set of features to detect new types of fake news. In summary, the adaptive approach should add and combine existing factors and provide explanations for what exactly influences the result of detecting fake news.

3. Methods and materials

At work, a new approach is proposed for detecting fake news, which is based on analyzing generalized characteristics of the content rather than just the text itself. To detect fake news, experts use a set of generalized content characteristics. Typically, the text is examined for faulty reasoning (arguments are supported by "rotten" evidence, quotes are attributed to unknown sources, numerical figures are presented without indicating their sources, etc.). Indicators of faulty reasoning include: theses not supported by credible evidence, common myths instead of arguments, lack of specific data and sources, and so on. Text can also be evaluated for emotionally charged content that manipulates the reader, with the goal of making the reader a "useful idiot." This is achieved through exaggeration, epithets, negatively connotated words, and strong emotional appeals that shut down the reader's logic and encourage them to act based on outrage. The industry of creating fake news is constantly evolving, and other methods of creating them are possible. Therefore, there is a need to propose an approach that would allow an expert to analyze the text based on its existing characteristics and also provide tools to add new characteristics and "retrain" classifiers on new sets of fake news.

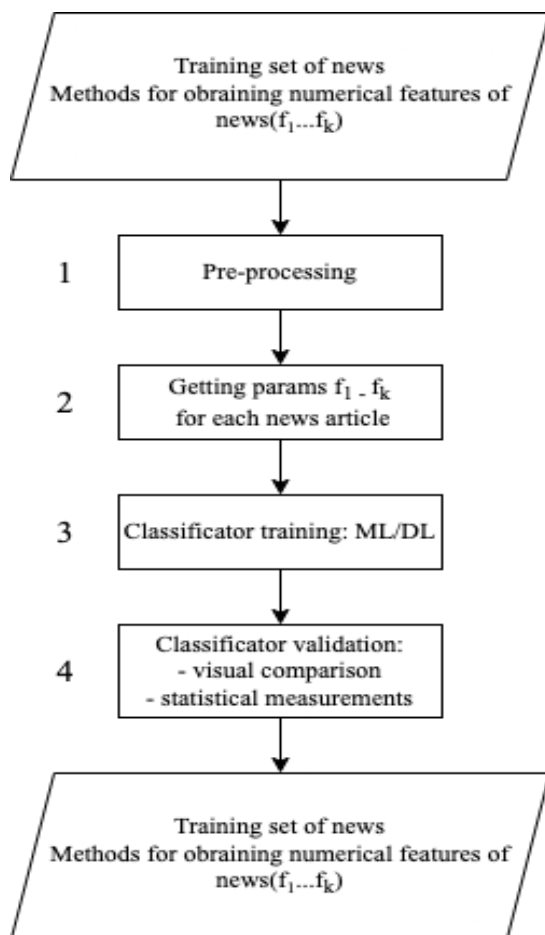


Figure 1: Scheme of the classifier training method

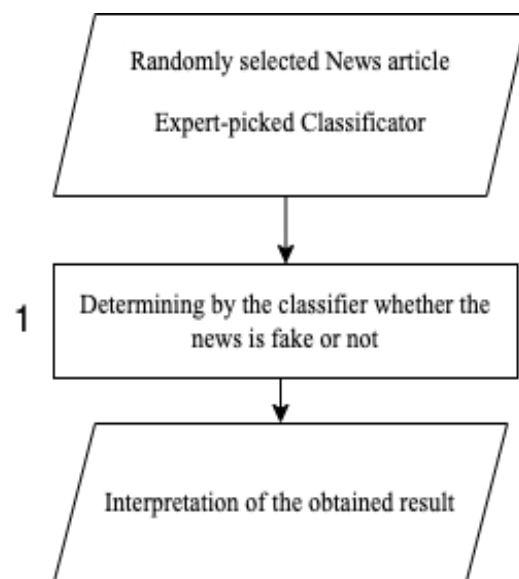


Figure 2: Scheme of the classification method

The proposed approach consists of a method of training classifiers (based on various characteristics of the text and training data sets) (Figure 1) and a method of classification using the selected classifier (Figure 2).

As can be seen from Figure 1, the input information for the classifier training method consists of a training set of labeled fake and non-fake news and a set of methods for obtaining numerical characteristics of the text. The next step is text preprocessing. Then, the news text is transformed into a feature vector using the methods of obtaining numerical characteristics. The resulting labeled set of feature vectors is fed into the classifier. The classifier can be any machine learning (ML) or deep learning (DL) method. The resulting classifier is analyzed for its ability to classify both the training and testing data sets. After evaluation, the classifier can be used for detecting fake news.

For the classification method (Figure 2), the input information is arbitrary news text and a classifier selected by an expert. The result of the method is to determine whether the news belongs to a fake or non-fake category.

Further, we will describe in detail the main steps of the presented methods and the algorithms and methods used in the research for transforming text information.

3.1. Textual content analysis tools

To analyze the ability of the proposed approach to detect fake news, the spaCy Python NLP library [8] was used, which includes a range of natural language processing tools, including named entity recognition, part-of-speech tagging, and dependency parsing. The large spaCy English model was used, which includes pre-trained word embeddings that can be used for computing similarity between texts, as well as the spacytextblob [9] library for determining sentiment and polarity. Additionally, scikit-learn [10] was used for computing Multidimensional Scaling [11] and Support-Vector Machine [12]. Although there are several NLP libraries available, the use of spaCy and scikit-learn was due to their ease of use and access to pre-trained models, such as the BERT base model, and the ability to work with a pre-trained Ukrainian model. Other alternative libraries could include NLTK, Stanford CoreNLP, and Gensim. However, the analysis showed that spaCy provides the best combination of performance and ease of use to achieve the research goal.

3.2. Pre-processing

The first step in preparing text for NLP processing involves cleaning the text and removing any irrelevant or unnecessary information [13]. This typically involves removing punctuation marks, numbers, and stop words, which are common words that do not carry much meaning, such as "the," "and," and "a." In the proposed approach, the built-in stop word list from spaCy is used to remove stop words from the text. Removing stop words is important because it can help reduce noise in the text and facilitate the identification of important words and phrases. After cleaning the text, it is tokenized using the spaCy tokenizer, which breaks the text into individual tokens or words. Each token is assigned a part-of-speech tag that indicates the role the word plays in the sentence, such as noun, verb, or adjective. Next, the spaCy lemmatizer is used to reduce each token to its base form or lemma. Lemmatization is important because it can help reduce the complexity of the text and facilitate comparisons between words that have the same root or meaning. These processing steps can be useful in detecting fake news by facilitating the identification of important words and phrases in the text and removing irrelevant or unnecessary information. Reducing the complexity of the text and identifying key words allows for better detection of patterns and features in the text that indicate fake news or biased language.

3.3. Characteristics of textual content

Next, the set of text characteristics used in this study will be considered. It should be noted that it is not fixed. These characteristics are used to analyze the ability of the proposed approach to solve the task at hand. It should also be noted that the proposed approach is adaptive, allowing for the expansion of both the set of text characteristics and the training data sets.

3.3.1. Persuasion and influence

Fake news can be convincing and influential because they often use language and tactics aimed at manipulating the reader's emotions and beliefs. For example, fake news can use biased language to appeal to the reader's existing beliefs and values, or use persuasion techniques such as repetition, paraphrasing, and dehumanizing language to influence the reader's perception of the topic.

Language bias refers to language that expresses preference or bias towards a particular group or belief system. In fake news, biased language can be used to draw attention to readers who share similar beliefs or values, as well as to reinforce the beliefs of those who already agree with the message. For example, a news article that criticizes a particular political figure may use derogatory language to appeal to readers who already oppose that figure, while also strengthening negative beliefs among readers.

Subjectivity is another important factor in biased language, as it can complicate an objective evaluation of the content of a news article. Fake news can use intentionally subjective or emotional language to sway the reader's opinion or beliefs. For example, an article that presents a certain political figure in a negative light may use language intended to provoke the reader's feelings of anger or sadness in order to influence their beliefs about that figure.

Other methods commonly used in fake news include paraphrasing, repetitive narratives, dehumanizing language, and objectification. These techniques can be used to reinforce the message of an article and make it more memorable and influential for the reader. For example, an article that criticizes a certain group may use dehumanizing language to make the group seem less sympathetic or relatable, making it easier for the reader to dismiss their concerns or opinions.

According to the given characteristics of the text, it is suggested to use the following parameters:

- f_1 - `paraphrased_ratio`: the paraphrasing coefficient allows you to find the percentage of information that has already been voiced, but is repeated for some purpose; this parameter was calculated by comparing the previous sentence with the following sentences; measured from 0 to 1, where 0 is no paraphrasing and 1 is a complete repetition of the text;
- f_2 - `dehumanizing_language_ratio`: coefficient of "deprivation" of human contact; this parameter was a computational measurement of the Proper Noun (grammatical construction) in the sentence and the mismatch of the Part of speech; is measured from 0 to 1, where 0 is normal handling and 1 is maximum dumanization;
- f_3 - `subjective_words_ratio`: the coefficient of subjective words shows the subjectivity of the text; determined using the `spacytextblob` component, which contains a ready-made subjectivity indicator for English words; measured from 0 to 1 according to increasing subjectivity.

3.3.2. Narrative

Narrative is one of the components of detecting fake news. It is important for news to have a clear and consistent narrative that is related to the headline and the overall essence of the text. The narrative is revealed through analyzing the context of the news and describes the logical order of events or information contained in the text.

Special attention should be paid to the narrative in the case of fake news, as they may contain illogical connections between the information that reaches the reader and the headline. Fake news often contains attempts to change the audience's opinion or create a nonexistent problem, which can lead to social division or panic. In such cases, the narrative may be inconsistent and illogical, which is a sign of a fake.

In analyzing the narrative, it is important to evaluate not only the connections between the news headline and the text, but also the logical connections between events and facts presented in the text. This makes it possible to detect fake news that may contain illogical and conflicting connections between facts and events.

Therefore, detecting fake news depends on how clearly they are structured and logically connected. The more attention is paid to the narrative, the greater the possibility of detecting fake news and preventing the spread of false information among the audience.

According to the given characteristics of the text, it is suggested to use the following parameter:

- f_4 - header_summary_similarity_ratio: the similarity coefficient of the title of the article and its body - determined by comparing the title and body of the article using the similar method of the spaCy library; generalization of the body of the article is determined by the selection of the most important sentences based on their similarity to the rest of the article; is measured from 0 to 1, where 0 is dissimilarity and 1 is identity of title and body.

3.3.3. Sentiment and Linguistic Analysis

Sentiment analysis and linguistic analysis are widely used methods in detecting fake news. An important part of these methods is identifying unusual and illogical textual structures, as fake news may contain vague and unmotivated statements that contradict the headline or general idea of the news.

To analyze the sentiment of fake news texts, methods that allow the determination of the average, positive, negative, and neutral mood are used. Machine learning algorithms and language analysis are typically applied for this purpose. Identifying such parameters helps to distinguish fake news from real news because fake news may have an overly positive or negative sentiment that does not correspond to the content of the news. These methods are an important tool in combating the harmful effects of fake news on society and enable informed conclusions to be made about the veracity of the text.

According to the given characteristics of the text, it is suggested to use the following parameters:

- f_5 - unusual_inappropriate_language_ratio: the coefficient of unusual inappropriate language shows how many unusual words there are; determined by checking tokens (words and not only) to see if they fall under the standard category: is_alpha, is_punct, exists in vocabulary; measured from 0 to 1 according to the number of words in the entire text;
- f_6 - awkward_text_ratio: the coefficient of awkward, complex or convoluted sentence structures is determined by taking into account and subtracting the dependencies of linguistic tagging in the text "amod", "compound", "nsubj", "dobj", and "pobj"; measured from 0 to 1 according to the number of complex "tokens";
- f_7 - avg_sentiment: the sentiment coefficient of the text of the words shows the average sentiment of the text; determined using the spacytextblob component, which contains a ready-made polarity indicator; is measured from -1 to 1, where -1 is negative, 0 is neutral, and 1 is positive;
- f_8 - positive_ratio: the positivity coefficient shows how positive the text is; determined by subtracting the number of positive words from the entire text; is measured from 0 to 1;
- f_9 - neutral_ratio: the neutrality coefficient shows how neutral the text is; determined by subtracting the number of positive words from the entire text; is measured from 0 to 1;
- f_{10} - negative_ratio: the negativity coefficient shows how negative the text is; determined by subtracting the number of positive words from the entire text; is measured from 0 to 1.

3.4. Evaluation of the validity of the proposed feature vector

To assess the quality of the proposed feature vector for classification tasks, the use of Multidimensional Scaling (MDS) method is proposed. This is one of the methods for reducing the dimensionality of the vector space. The aim of the method is to reduce the dimensionality to a level that can be visualized (3D or 2D). The criterion for dimensionality reduction is, for example, the Euclidean distance between vectors. That is, by solving an optimization problem, an $R^n \rightarrow R^2$ mapping is found that makes it possible to obtain a two-dimensional graph of the mutual arrangement of vector points and visually assess the quality of the model for the classification task. Visual criteria have been proposed to assess the quality of modeling (Figure 3).

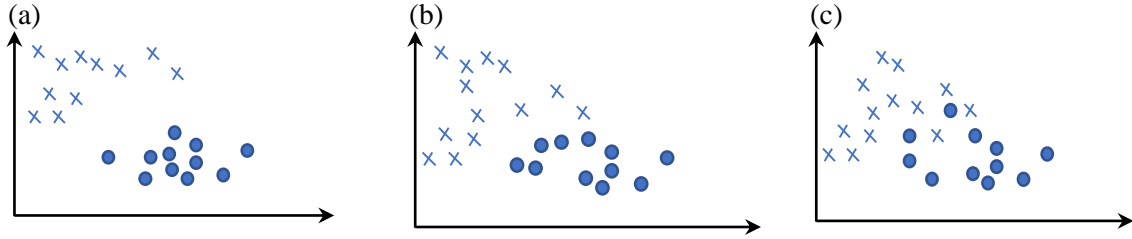


Figure 3: The quality of the feature vector for the classification problem (a) is ideal, (b) is acceptable, (c) is satisfactory

Criteria 1. – An ideal feature vector for text classification. Figure 3 (a) shows that the two classes are clearly separated.

Criteria 2. – Acceptable feature vector for text classification. Figure 3 (b) shows that two classes collide with each other, but individual members of the classes do not intersect.

Criteria 3. – Satisfactory model level for text classification. Figure 3 (c) shows that the two classes overlap somewhat. With such an indicator, the model can be considered workable, but it will require an additional expert opinion to confirm the classification.

The proposed criteria are recommended to be used to verify the quality of the proposed feature vector. The feature vector will be considered correct if the values of the results appear as shown.

Subsequently, if the feature vector allows for the separation of two classes of news, a classifier is proposed to be obtained.

The next step is to evaluate the quality of the proposed classifier using the following metrics: precision, recall, and F_1 -norm.

In machine learning precision and recall are indicators of productivity [14]. They apply to information obtained by simple sampling, collection or corpus.

Precision shows what proportion of the results found from the sample are relevant to the query [15], and is by the formula:

$$precision = \frac{relevant\ document \cap\ retrieved\ document}{retrieved\ documents} \quad (2)$$

The best result for classification issues is a score of 1.0, when each of the samples submitted for entry actually belongs to a certain class (however, the number of such samples that was not observed correctly is unknown). Relevant documents can still be called correctly classified.

Recall shows the share of relevant documents that found are successfully [16], and is formally depicted as follows:

$$recall = \frac{relevant\ document \cap\ retrieved\ document}{relevant\ document} \quad (3)$$

The F -measure is calculated through precision and recall. It is common to use the measure F_β , in which β depending on its value, pays more attention to either precision or recall. However, they often focus on the measure F_1 . Measure F_1 – is the weighted harmonic mean of precision and recall [17] which can be formally written as follows:

$$F_1 = 2 \times \frac{recall \times precision}{recall + precision} \quad (4)$$

The best score for F_1 is 1.0, which suggests that precision and recall are ideal. The worst score is 0 if either precision or recall is zero. Given the popularity of measure F_1 , it should be noted that it can give inaccurate data with an unbalanced data set, so it should be used only on a balanced set [18].

These metrics are used to study the results.

4. Results and discussion

A number of experiments were conducted to test the proposed approach and evaluate the validity of the feature vector. Below are their results and discussion. A description of the dataset used in the experiments is given. The result of the application of visual analytics to assess the ability of the proposed features of fake news texts to be divided into two classes is given. Visual and numerical (statistical metrics) results of classifier training (using SVM) are given. The discussion was carried out and the prospects of the proposed approach were given.

4.1. Dataset

The dataset[19] has over 20000 true and fake news labeled and categorized. It is very popular among the data science community and has been used in many articles and works.

4.2. MDS

The results of applying the MDS method to input data (generalized features of news texts) in 2-dimensional space are shown in Figure 4.

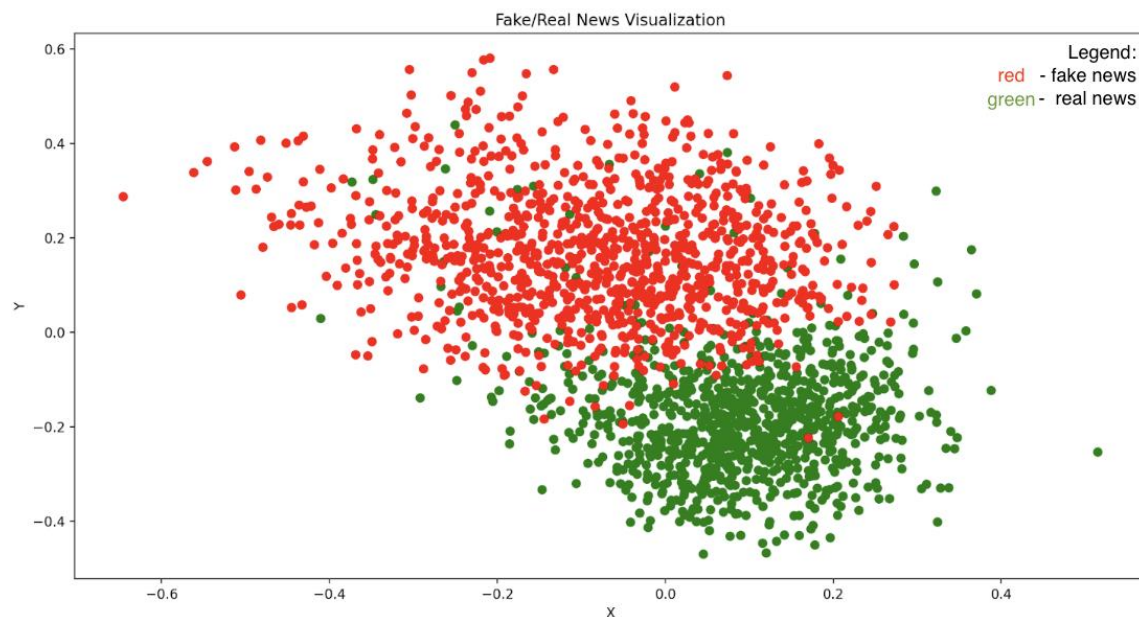


Figure 4: MDS Results for 2000 articles

As can be seen from Figure 4, the result is satisfactory, the classification was successful for a larger number of texts from the training set. Analysis of a small number of misclassified texts showed that there are true articles written with poorer text quality and vice versa.

4.3. SVM

After calculating the MDS, we can pass a value to the `train_test_split` method to split the data into training and test samples. Using SVM methods from the `scikit-learn` library, we obtained the following results:

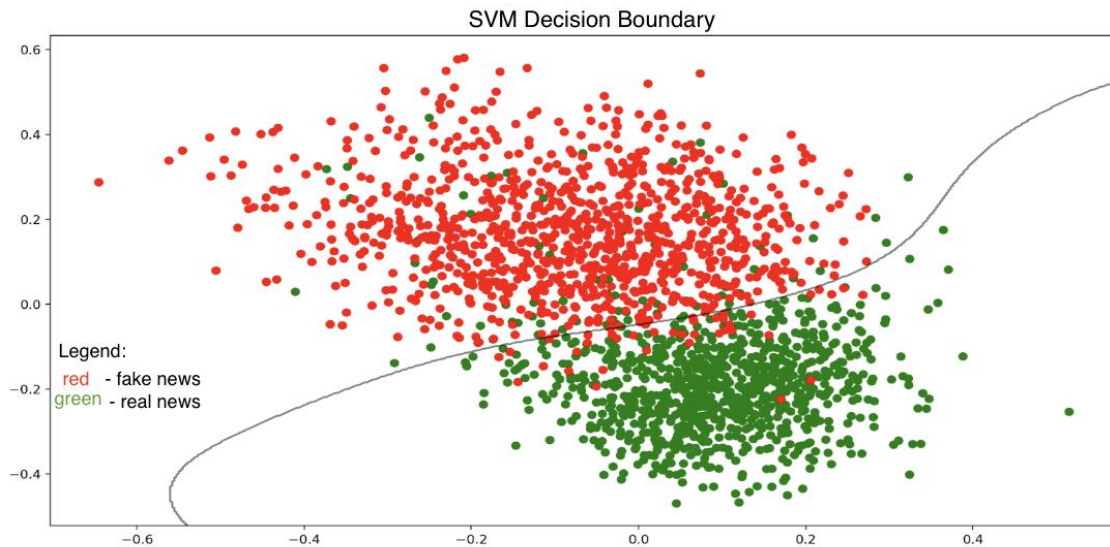
Table 1

Comparison of the metrics for the classification problem

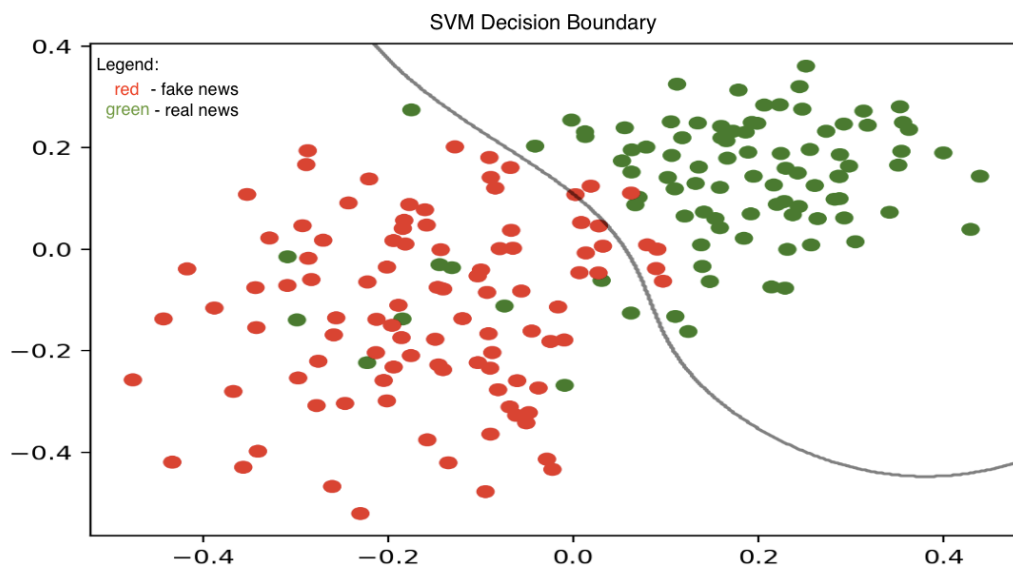
| Number news, N | Precision | Recall | F_1 |
|----------------|-----------|--------|-------|
| 20 | 1.0 | 1.0 | 1.0 |
| 200 | 0.88 | 0.82 | 0.85 |
| 2000 | 0.93 | 0.92 | 0.93 |

After the number of news articles went over 2000, the results became consistent and we could consider it average for the whole dataset.

The obtained numerical results show the high accuracy of the proposed approach for determining fake news. The given values of the statistical metrics are either in the range or even better than the published modern results of other researchers.

**Figure 5:** SVM decision boundary for 2000 elements

Also, in Figure 5 and 6, we can outline the decision boundary [20]. The boundary is determined by the support vectors, which are the data points closest to the hyperplane. The SVM then uses this decision boundary to classify new, unlabeled data points based on which side of the boundary they fall on.

**Figure 6:** SVM decision boundary for 200 elements

4.4. Limitations of the approach and further research

The main limitation of the proposed approach is the lack of high-quality labeled datasets for successful training of classifiers, especially for the Ukrainian language. Another limitation is the insufficient number of generalized characteristics of texts that allow detecting more hidden ways of creating fakes. However, it should be noted that these limitations are not significant for the proposed approach since it allows for adaptation, building new interpreted and transparent classifiers using both new datasets and additional generalized text features.

The future development of the approach to fact-checking may include the integration of external APIs to gather more detailed information and fact-check claims made in articles. These APIs may be from verified sources such as news agencies, government organizations, or other fact-checking organizations. This will help improve the accuracy and reliability of the fact-checking process.

Another potential development could be checking information on different social media platforms such as Twitter to verify the popularity and authenticity of claims. This can be done by analyzing the number of likes, retweets, and article publications, as well as verifying the sources of information to ensure their reliability. Additionally, the approach can also detect the toxicity of comments on social media platforms such as Twitter.

Finally, the approach can be extended to detect content created by artificial intelligence (AI). This may involve analyzing the language and structure of text to detect patterns that are commonly used in AI-generated content. Detecting AI-generated content will help prevent the spread of misinformation and disinformation.

5. Conclusion

An adaptive approach to identifying fake news using natural language processing techniques and machine learning algorithms is presented in this work. A thorough review of related works was conducted to ensure the novelty and effectiveness of the proposed approach. Ten different parameters (general text features) were used to model the text, and multidimensional scaling (MDS) was applied to obtain visual analytics as one of the criteria for evaluating the quality of the proposed approach. A support vector machine (SVM) classifier was trained to classify text into different categories. The research results show that the proposed approach is within the same range or surpasses existing methods in accuracy (overall accuracy - over 90%).

Limitations of the proposed approach include the absence of high-quality, annotated datasets (especially for the Ukrainian language) for successful classifier training and insufficient generalized text features (for detecting more hidden ways of creating fake news). These limitations are not critical as the proposed adaptive approach is capable of incorporating new datasets and new generalized features for retraining.

Future improvements to the approach will be directed towards increasing the accuracy of identifying fake news and achieving greater interpretability and understanding of classification results.

6. References

- [1] Haoran Wang, Yingdong Dou, Canyu Chen, Lichao Sun, Philip S. Yu, Kai Shu, Attacking Fake News Detectors via Manipulating News Social Engagement (2023). <https://doi.org/10.48550/arXiv.2302.07363>
- [2] Alim Al Ayub Ahmed, Ayman Aljabouh, Praveen Kumar Donepudi, Myung Suh Choi: Detecting Fake News Using Machine Learning (2021). <https://doi.org/10.48550/arXiv.2102.04458>
- [3] Shaina Raza & Chen Ding: Fake news detection based on news content and social contexts: a transformer-based approach (2022): <https://doi.org/10.1007/s41060-021-00302-z>
- [4] Federico Monti, Fabrizio Frasca, Davide Eynard, Damon Mannion, Michael M. Bronstein: Fake News Detection on Social Media using Geometric Deep Learning, <https://doi.org/10.48550/arXiv.1902.06673>

- [5] Mayank Kumar Jain; Dinesh Gopalani; Yogesh Kumar Meena; Rajesh Kumar: Machine Learning based Fake News Detection using linguistic features and word vector features, <https://ieeexplore.ieee.org/document/9376576>
- [6] Hnin Ei Wynne , Zar Zar Wint: Content Based Fake News Detection Using N-Gram Models: <https://dl.acm.org/doi/10.1145/3366030.3366116>
- [7] Xinyi Zhou , Atishay Jain , Vir V. Phoha , Reza Zafarani, Fake News Early Detection: A Theory-driven Model, <https://dl.acm.org/doi/10.1145/3377478>
- [8] spaCy, Python library for NLP processing , <https://spacy.io>
- [9] spacytextblob, Sentiment analysis component for spaCy, <https://spacy.io/universe/project/spacy-textblob>
- [10] Ski-learn, classification and other library for SVM, MDS <https://scikit-learn.org/stable>
- [11] MDS(Multidimensional Scaling), Mead, A (1992). "Review of the Development of Multidimensional Scaling Methods". Journal of the Royal Statistical Society. Series D (The Statistician). 41 (1): 27–39. [JSTOR 234863](https://www.jstor.org/stable/234863).
- [12] SVM (Support Vector Machine), [Cortes, Corinna](#); [Vapnik, Vladimir](#) (1995). "Support-vector networks" (PDF). *Machine Learning*. 20 (3): 273–297. [CiteSeerX 10.1.1.15.9362](#). [doi:10.1007/BF00994018](https://doi.org/10.1007/BF00994018). [S2CID 206787478](#).
- [13] Text Preprocessing in Python using spaCy, <https://iq.opengenus.org/text-preprocessing-in-spacy>
- [14] R. Yacouby, D. Axman, Probabilistic Extension of Precision, Recall, and F₁ Score for More Thorough Evaluation of Classification Models, Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems, Association for Computational Linguistics (ACL), Stroudsburg, PA, USA, 2020, pp. 79-91. doi:10.18653/v1/2020.eval4nlp-1.9.
- [15] R. Padilla, S. L. Netto, E. A. B. da Silva, A Survey on Performance Metrics for Object-Detection Algorithms, 2020 International Conference on Systems, Signals and Image Processing (IWSSIP), 2020. doi:10.1109/IWSSIP48289.2020.
- [16] J. Miao, W. Zhu, Precision-Recall Curve (PRC) Classification Trees, Evolutionary Intelligence, 2021. 10.1007/s12065-021-00565-2.
- [17] R. Aliguliyev, R. Aliguliyev, Ya. Imamverdiyev, L. Sukhostat, An improved ensemble approach for dos attacks detection, Radio Electronics, Informatics, Management 2: 2018, pp. 73–82. doi: 10.15588/1607-3274-2018-2-8.
- [18] A. Tharwat, Classification assessment methods, Applied Computing and Informatics Vol. 17 No. 1, 2021, pp. 168-192. doi:10.1016/j.aci.2018.08.003.
- [19] Clement Bisailon, Fake and real news dataset, 2019, <https://www.kaggle.com/datasets/clmentbisailon/fake-and-real-news-dataset>
- [20] SVM Decision Boundary, https://scikit-learn.org/0.18/auto_examples/svm/plot_iris.html