

Comparative Analysis of the Efficiency of Techniques for Detecting Misinformation in Healthcare Data*

Alžbeta Žiarovská

Slovak University of Technology in Bratislava

Faculty of Informatics and Information Technologies

`xziarovska@stuba.sk`

October 31, 2023

*Semestral project in subject Engineering Methods, ac. year 2023/24, guidance: MSc. Mirwais Ahmadzai

Abstract

...

1 Introduction

In this article I discuss the current situation regarding spread of misinformation in the medical field. This topic is very important in the aftermath of the global COVID-19 pandemic, as we have seen a great rise of various misinformation on the internet, which provide danger to society or even lives [8]. The main problem in my perception is, that the easy access to all the information on the Internet, which does not necessarily has to be true, can increase fear and anxiety and ultimately lead to the delay of diagnosis and receiving the effective healthcare in the case the information are not perceived correctly [7].

2 Related work

The topic of misinformation is quite often researched. Some of the work is focused on comparing the different machine learning techniques [6], [5]. Focus of these researches are to introduce how these methods work. Some articles bring new way of misinformation recognition technique [2] by combining the two, or even more techniques. Plenty of research have been also done in a field of medical misinformation, which is a important part of my article [4], [3], [7]. These articles mostly focus on defining misinformation their possible cause and consequences.

3 Methodology

I am focusing on analyzing machine learning models as a way to find the medical misinformation. The two machine learning techniques I am focusing on are Naïve Bayes and Support Vector Machine. I introduce these techniques and compare their efficiency in order to establish which one is the most suitable for healthcare misinformation detection. The comparison is done by comparing the accuracy rates of detecting false information among true information.

4 Misinformation in healthcare

Difference between misinformation and disinformation The terms *misinformation* and *disinformation* are much the same, however, a small, but crucial difference can be distinguished. The difference between the two is a intention with which the false information is made accessible to the public and spread. Whilst the misinformation is usually created without direct intention of misleading and spreading false, meaning the person who put the information into the world might not actually know it is not true. On the other hand, disinformation is essentially created to spread false information. An example of such activity can be political propaganda [4] [3]. Even though the terms are not meaning the same, for the purpose of this article they are used as synonyms, because the author's knowledge, whether the information is factual, is negligible in the scope of its false recognition.

Health care misinformation A vast majority of people is using the Internet and social media for entertainment or information seeking. However, with the possibility of immediate communication and sharing, it has become easy to spread misinformation online [7]. During the COVID-19 pandemic there have been a great amount of healthcare misinformation spread regarding vaccines and their effectiveness [2]. Internet is easily accessible and more and more people are looking for relevant health information without the proper knowledge of how to distinguish, whether the information is true. This can lead to unintentionally getting false information, as many websites do not provide accurate medical information [3]. Another example of current situation can be the popular misinformation about vaccines causing autism, which was repeatedly proven as nonfactual information [7]. The spread of medical misinformation is not only occurring in the 21st century. In the past there was false information about public health impact of smoking spread by tobacco companies, which was later proven as false. [3].

5 Misinformation recognition techniques

5.1 Machine learning techniques

5.1.1 Naïve Bayes

The Naïve Bayes method is a linear probabilistic machine learning technique based on Bayes theorem. This method uses probability of the events without taking their relation into the consideration [6]. This approach might not look to be the best, as the words have their order and are related one to other in articles. However, the opposite is true, as the linear models are capable of achieving high efficiency despite their simplicity [5]. The accuracy of Naïve Bayes (as well as other machine learning methods) is also depended from which type of measuring the importance of the words in the documents is used. For the Naïve Bayes the results vary from 84,056% [6] to 98,71% [1]. These percentages represent the accuracy of distinguishing false and true information by machine. The closer analysis of these differences and their comparison is given in Section ??.

Formula for Naïve Bayes calculation: [6]

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1)$$

$P(A|B)$ is the probability of event A happening supposing, that event B has occurred.

5.1.2 Support Vector Machine

Support vector machine might be classified as a binary technique, as its methodology is to divide the data it was given into two categories [5] (in the case of misinformation detection into true and false information). The division is made by creating a hyperplane (a object in the vector space with one dimension less, that the vector space itself [6]) As it was mentioned in the section 5.1.1 about Naïve Bayes, the result can vary according to the technique used for analysis of the given data

and for the Support Vector Machine percentages of accuracy are on a scale from 83% [2] to 95,05% [6].

6 Analysis and results

I have introduced a couple of machine learning techniques to detect and highlight misinformation in healthcare in order to compare their effectiveness. To do so, I have searched of numeric expression of their accuracy in percentage - how big is the proportion of correctly categorized claims or documents. In Chaphekar's [2] and Poddar's [5] there have been used two different methods to establish the importance of words in document. That is the reason why there are two different results for both of the machine learning techniques in these papers. I have made an average of these percentage, which you can find in Table 1 with corresponding superscripts.

Naïve Bayes	Support Vector Machine
88.37% ¹	84% ¹
98.71% ²	94.17% ²
85.85% ³	90.95% ³
84.06% ⁴	95.05% ⁴

Table 1: Accuracy of machine learning techniques in misinformation detection according to various researches, 1 - [2], 2 - [1], 3 - [5], 4 - [6]

According to the results of the table, Support Vector Machine is more efficient in detecting misinformation with average accuracy rate of 91.04%. However, the Naïve Bayes method is behind just by a little bit the average percentage of 89.25%. These result are not so greatly varying from each other and it shows, that machine learning techniques nowadays are highly advanced and can achieve brilliant results.

7 Discussion

I have found several evidence, that machine learning techniques can be a very useful (and accurate) tool in recognizing misinformation in healthcare. The impacts of spreading false medical information can be devastating and it is important to spread awareness of their existence and how to effectively recognize them and assume the correct approach. The world is changing quickly and what is true today might not be correct tomorrow and technology mentioned in this article provides a serious help in this direction. There are already many websites available using these techniques to provide a relevant output to help us find the true information on the medical question we have.

8 Conclusion

This article is a literature review of articles researching misinformation detection in healthcare data. I have compared results from several resources and based on these I have created an analysis of two machine learning techniques I focused on. In order to provide analysis of their efficiency I came to conclusion, that the machine learning techniques are highly advanced and can differentiate between factual and nonfactual information efficiently. This can provide solution to the main problem I mentioned in the beginning of the article. People may not always perceive correctly the information found on the Internet.

References

- [1] Yashoda Barve and Jatinderkumar R Saini. Healthcare misinformation detection and fact-checking: a novel approach. *International Journal of Advanced Computer Science and Applications*, 12(10), 2021.
- [2] Garima Chaphekar. *Unmasking Medical Fake News Using Machine Learning Techniques*. PhD thesis, San Jose State University, 2022.
- [3] John Cook, Ullrich Ecker, and Stephan Lewandowsky. Misinformation and how to correct it. *Emerging trends in the social and behavioral sciences: An interdisciplinary, searchable, and linkable resource*, pages 1–17, 2015.
- [4] Andrew M Guess and Benjamin A Lyons. Misinformation, disinformation, and online propaganda. *Social media and democracy: The state of the field, prospects for reform*, 10, 2020.
- [5] Karishnu Poddar, KS Umadevi, et al. Comparison of various machine learning models for accurate detection of fake news. In *2019 Innovations in Power and Advanced Computing Technologies (i-PACT)*, volume 1, pages 1–5. IEEE, 2019.
- [6] Jasmine Shaikh and Rupali Patil. Fake news detection using machine learning. In *2020 IEEE International Symposium on Sustainable Energy, Signal Processing and Cyber Security (iSSSC)*, pages 1–5. IEEE, 2020.
- [7] Yuxi Wang, Martin McKee, Aleksandra Torbica, and David Stuckler. Systematic literature review on the spread of health-related misinformation on social media. *Social science & medicine*, 240:112552, 2019.
- [8] Haider Warraich. Dr. google is a liar. *New York Times*, 2018.