


# TA-WHI: Text Analysis of Web-Based Health Information

Piyush Bagla, Dr. B.R. Ambedkar National Institute of Technology, Jalandhar, India\*

 <https://orcid.org/0000-0002-1664-7787>

Kuldeep Kumar, National Institute of Technology Kurukshetra, India

## ABSTRACT

The healthcare data available on social media has exploded in recent years. The cures and treatments suggested by non-medical experts can lead to more damage than expected. Assuring the credibility of the information conveyed is an enormous challenge. This study aims to categorize the credibility of online health information into multiple classes. This paper proposes a model named Text Analysis of Web-based Health Information (TA-WHI), based on an algorithm designed for this. It categorizes health-related social media feeds into five categories: sufficient, fabricated, meaningful, advertisement, and misleading. The authors have created their own labeled dataset for this model. For data cleaning, they have designed a dictionary having nouns, adverbs, adjectives, negative words, positive words, and medical terms named MeDF. Using polarity and conditional procedure, the data is ranked and classified into multiple classes. The authors evaluate the performance of the model using deep-learning classifiers such as CNN, LSTM, and CatBoost. The suggested model has attained an accuracy of 98% with CatBoost.

## KEYWORDS

Credibility, Data Mining, Health information, Machine Learning, NLP, Online Health, Text Mining, WHI

## INTRODUCTION

In the healthcare industry, wrong treatment, misinformation, self-treatment, and myths related to unconventional treatments is not a recent development. It is as ancient as medical care itself. Before the boom of the Internet, Radio, and Television, this issue was based on the therapeutic relationship as well as its context (Fernández-Celemín & Jung, 2006). The spectrum of damage is taken to an entirely new degree because of global technological advancement. Misinformation on social media became so common that in 2016 Oxford dictionary introduced “post-truth,” meaning “relating to or denoting circumstances in which objective facts are less influential in shaping public opinion than appeals to emotion and personal belief” (Harsin, 2018). Posting misleading or misinformation on social media is a fashion for some.

Social networking services like Facebook, followed by Twitter, are currently the industry leaders, with over 1.3 billion members and a monthly average fluctuation of 300 million people. Every second,

DOI: 10.4018/IJSSCI.316972

\*Corresponding Author

This article published as an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

their interactions create gigabytes of data (Alrubaian et al., 2018; Ranganath et al., 2017). Online social networks are appealing because they provide a quick and easy way to acquire health information. It is also quite simple to share information with others. However, the broad dissemination of incorrect information is made possible by rapid data scattering at a high pace with little effort. Thanks to the pandemic in 2020, social media usage increased by many folds. More information is now shared on social media than before 2020 (Zhang et al., 2017). The world has seen how misinformation about COVID spreads like wildfire, and every time World Health Organization (WHO) or some medical authority comes up to deny the news. People are so scared to visit hospitals that they prefer the social media Doctor (Zwolenski & Weatherill, 2014). Following the incorrect therapeutic advice given on social media might be fatal.

Text analysis is the practice of analyzing a vast amount of textual material to capture the key concept, trends, and hidden relationships. It is the process of transforming the unstructured text into a structured format to identify meaning patterns and new insights. Analysis of text is a crucial step in getting the hidden meaning behind it. The most popular technique for doing so is sentiment analysis. There are a number of researchers who have used this technique to get the actual sentiment behind the post on social media, especially on Twitter. It is extended further to incorporate a machine learning algorithm to perform the classification task. As a result, the credibility of the post can be identified (Alharbi & Alhalabi, 2020; Gunti et al., 2022; Mohammed et al., 2022). People are now mixing their regional language while making any social media post, no matter from which country they belong. In technical terms, we call this code-mixing. This makes the analysis of text even more difficult. However, techniques such as Machine Learning, Neural Networks, and LSTM (Long Short-Term Memory) can be used to mitigate the problem of code-mixing (Sharma et al., 2021; Singh & Sachan, 2021). With the development of technology, the amount of data generated on the Internet has increased daily. This data includes valuable patterns that must be recognized to get meaningful information. There are several methods that facilitate the completion of this task, such as using data mining techniques (García-Peñalvo et al., 2021), text mining and privacy preservation techniques in name analysis (Veluru et al., 2015), and scientific issue tracking with topic analysis based on crowdsourcing (Kim et al., 2018). In one way or another, all these methods contribute to the data mining process. However, there is always room for strengthening the capabilities of the proposed approaches. For example, there is always a question regarding the authenticity of the information patterns found during text mining. Very few studies address this issue and those that do have several drawbacks. To determine the veracity of web-based health information, they used a predetermined data mining algorithm that operates on the existing dataset. However, there is an ongoing need to develop a strategy based on a user-generated algorithmic approach to determining the credibility of web-based health information using real-time datasets.

We may break up research on automatic reliability evaluation for online news into three sections. In knowledge bases, some people try to extract and validate allegations stated in the text. It is further known as “fact-check.” Others focus on measuring trustworthiness in the news context content’s social media environment and its creator (Atanasova et al., 2019; Ciampaglia, 2015; Thorne et al., 2018). The third option depends on the text’s general linguistic characteristics, such as writing style. Low-credibility material, such as fake news, is written to evoke an intense emotional response (Tacchini et al., 2017; Zubiaga et al., 2016), which requires specific stylistic methods. Measurement of linguistic complexity, detection of syntactic patterns using n-grams of part-of-speech tags or counting words belonging to particular categories might act as credibility indicators (Ahmed et al., 2021; Bhattarai et al., 2021; Potthast et al., 2017; Reis et al., 2019).

The remaining sections are organized as follows: First, we explore related work, including some of the most important research undertaken in this field. We then explain our proposed methodology, which is further broken into many sub-sections. The study’s results are presented in the next section, followed by the conclusion.

## RELATED WORK

Given the exponential rise of data in the health sector, a growing amount of effort has gone into developing automated procedures and strategies for analyzing data and extracting relevant information. Large data sets or databases are the focus of data mining; nevertheless, numerous specialized approaches have evolved, such as text mining and web mining, which are focused on text and web data, respectively (Zwolenski & Weatherill, 2014). Data mining techniques have been used in the healthcare domain to discover medical knowledge and patterns from clinical databases (Ling Liu & M. Tamer Özsu, 2018), text mining techniques are used to analyze unstructured data in the electronic health record, and web mining techniques have been used to investigate how people use healthcare-related websites and systems.

Researchers have explored using structured and unstructured data in the electronic health (Ahmed et al., 2021) record to understand and improve healthcare procedures. Some of the applications of data mining approaches for structured clinical data include the extraction of diagnostic criteria, identification of novel medical knowledge, and discovery of correlations between different types of clinical data. The researchers identified connections between operations done on a patient. They reported diagnoses using association rule generation, which may be beneficial for determining the efficacy of a set of procedures for diagnosing a specific condition. During public health emergencies, the endogenous health information demand created by internet users' lack of scientific understanding of health information promotes the distribution of health information by the media while also allowing rumor mongers to publish and propagate online rumors. G-SCNDR is an online rumor reversal model suggested by (Chen, 2016) to regulate the propagation of online rumors and decrease their harmful impact. Scientific knowledge level theory and an external online rumor control technique were used in the model. Several researchers have experimented with several techniques to detect fake news and rumors about Covid-19 and non-medical concerns transmitted through social media (Ahmed et al., 2021; Burel et al., 2021; Chawla et al., 2021; Kumari et al., 2021; Ling, 2021; Przybyła & Soto, 2021; Samadi et al., 2021; Wang et al., 2021; Yao et al., 2021; Zhou, Li, et al., 2021; Zhou, Xiu, et al., 2021). Since it is not the scope of the paper, we have just given the references.

Natural language processing (NLP) and text mining techniques have been used in healthcare for various purposes, including coding and billing, monitoring alternate courses of therapy, and identifying clinical problems and medical mistakes (Agle & Xiao, 2021). Several researchers have focused on developing text mining techniques for detecting certain types of co-occurring ideas in clinical documents (e.g., discharge summaries) and biological documents (e.g., illness medication or disease-finding). In a study, researchers discovered connections between illnesses and results (extracted from discharge summaries using a natural language processing tool). They utilized them to build a knowledge basis for an automated problem list summarizing system (Kulkarni Andrea, 2022). Another study uses free-text medical records to create illness profiles based on demographic data, primary diseases, and other clinical factors (Ehrenstein et al., 2022).

To determine if the feed is authentic or not, classification is important. Support Vector Machine (SVM) is very popular among researchers. SVM is used to classify the news (Saigal & Khanna, 2020). The favorite classifier of researchers, Convolution Neural Networks (CNN), converts the text into tokens and then classifies it. Researchers have frequently used CNN and hybrid models using CNN to classify fake news (Nasir et al., 2021; Yang et al., 2018). The researchers use transformer models like Bidirectional Encoder Representation (BERT) to extract the features from the text. The features are passed through the encoder layer and categorized the text using the dense layer (Kaliyar et al., 2021). Researchers used other transformer models like ALBERT and XLNET to classify fake COVID-19 news (Gundapu & Mamidi, 2021). Boosting algorithms have also been used for fake news detection and text classification. As the name suggests boosting algorithms boost a weak model. Gradient boosting (Xiong et al., 2018) effectively classifies fake news.

Many resources are available on the net to detect a fake feed (*FactCheck.Org - A Project of The Annenberg Public Policy Center, 2022; FACTLY - Making Public Data Meaningful, 2022; Full Fact, 2022; PolitiFact, 2022; Snopes.Com | The Definitive Fact-Checking Site and Reference Source for Urban Legends, Folklore, Myths, Rumors, and Misinformation., 2022; Vishvas News Fact Check Hindi: Fact Checking of Fake and Viral News, Photos & Videos, 2022; The Washington Post, 2022*). But they are not dynamic in nature, which means they are not updated immediately. Our study demands that academicians deal with the ensuing question: “For a complicated task like false web-based health news detection, which blend of pre-trained models and classifiers can work accurately?” Although some research studies have concentrated on deep neural methods, this field of study lacks comparison studies in health care that might pave the way for future research. This paper proposes a very simple yet very effective model to identify fake web-based health information.

## METHODOLOGY

There are the following objectives of the proposed work:

- Collect live data related to health.
- Apply NLP to retain relevant information.
- Classify the text as Neutral, Not Neutral, Misleading, Misinformation, and Good.
- Verify the classified data.

To achieve their objectives, we divided the methodology into three sections: Experimental setup, Objective function, and implementation.

### Experimental Setup

The setup includes hardware configuration, dataset, Tools, and Language.

#### Dataset

The datasets used to detect fake news are made up of real and fake news circulated across the Internet within a specified time. The subjects and content of datasets created at different times fluctuate because the general public’s interests heavily influence them. This causes issues with the model’s robustness. Because of the differences in word appearance, fake news detection models built from these datasets obtain high accuracy for datasets constructed for the same era and domain, but they lead to a reduction in the detection performance of false news in different domains and future applications (Murayama, 2021).

The dataset is collected from social media using the respective Application Programming Interface (APIs) of the social media applications. Relevant dictionaries are used to filter the related words. We collect the data dynamically to avoid minimizing the reliability of the data over time.

#### Tool and Language Used

To implement the proposed model, we used Python. For verifying the classification of the text, we used Buzzsumo.

Figure 1. TA-WHI model



### *Model*

The proposed model is summarized in the following steps:

- Collect feeds from the social media.
- Using NLP retains the text only.
- Filter each feed using Medical Data Framework (MeDF) dictionary.
- Retain words with frequency > 5.
- Calculate polarity.
- Rank the features.
- Classify using a different algorithm.
- Predict the new feed using the model designed.

### *Objective Function*

The objective function is represented mathematically:

$$\text{Class} = \left[ P(F_1, F_2, \dots, F_n) \Rightarrow R(P_1, P_2, \dots, P_n) \Rightarrow C(R_1, R_2, \dots, R_n) \right]$$

Here:

- $P$  is a polarity function that marks the polarity of each feed.
- $R$  is a Ranking function.
- $C$  is a classification function.
- $F$  is feed.

TA-WHI Filters the feed collected from social media then calculates the polarity and ranks them to get the class of the feed. In the objective function, feed  $F$  is defined as  $F = (e_1, e_2, e_3, \dots, e_n)$  where  $e$  is the extracted feature such that if word  $\in MeDF$  and frequency of  $w$  more than five, then store in  $e$ . where:

$$MeDF = \{noun, verb, pronoun, adjective, adverb, medical\ word\}$$

### *Steps Involved*

The work is summarized in the following steps:

- Step1:** Collect social media feed using respective APIs.
- Step2:** Remove special characters, numbers, and emoticons from the feed.
- Step3:** Split the sentences into words, using positive-negative adjective dictionaries and medical dictionaries, and retain the relevant word.
- Step4:** Calculate polarity.
- Step5:** Rank the feeds.
- Step6:** Finally, classify the feeds.



Here,  
M is medical word  
N is Neutral  
NN is not Neutral  
ML is misleading  
MI is Misinformation  
G is Good  
Step 8. Based on eq (2.1), the classes are assigned to the feed as:  
    if ML and post-truth on B, then  
news ← fabricated  
    else if ML and fake on B, then  
news ← fabricated  
    else if MI and misleading on B  
news ← Misleading  
    else if NN and advertisement on B, then  
news ← Advertisement  
    else if N and no negative Feed on B, then  
news ← Irrelevant  
    else if G and no negative Feed on B, then  
news ← Sufficient  
Step 9. Repeat steps 1 through 7 for all the features  
Step 10. Split the data created in 8 into a training set and a testing set  
Step 11. Test different classification algorithms on the Training set and testing set  
Step 12. Collect a new feed  
Step 13. Pre-process as in step 2  
Step 14. Run the classified model to predict the new feed for different algorithms

The text collected from the social networking sites like Twitter has irrelevant feeds like images, emoticons, special characters, and numbers. The feeds are to be cleaned such that our objective is to check the authenticity of the medical feeds. To achieve the objective, we created a dictionary MeDF. Using algorithm 1, we cleaned the text and then tagged and classified it:

retained\_words ∈ english\_dictionary

## RESULTS

Results include mathematical formulation and results based on them.

### The Mathematical Formulation for Accuracy

We calculated True Negative (TN), True Positive (TP), False Negative (FN), and False Positive (FP) to determine accuracy. The following equations are used to measure precision, recall, accuracy, and F1 score:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (4)$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (5)$$

$$\text{F1-Score} = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (6)$$

TP stands for the number of correctly evaluated authenticity. The number of sentiments detected by the model (FP) is not the same as that of genuine statements. FN denotes the number of statements that do not match. The output is a statement but not the authentic statement; TN is the number of erroneous authenticities discovered. F1-Score reflects the balance between precision and recall.

### Result Analysis

Twitter feeds were collected for three different diseases: COVID-19, Heart Attack, and Sugar. To run a machine learning algorithm, the data should be tagged correctly. Using steps 1-8 of Algorithm 1, the emotions were determined and classified. The tagged data is then split into training and testing sets.

The accuracy of TA-WHI relies on correctly identifying the tweets as negative, positive, and neutral. We used a threshold value to tag the feeds. To understand how accurately we marked the feeds as negative, positive, and neutral, we tested the feeds with BERT, Bio-BERT, and TA-WHI. We used pre-trained models BERT and BIO-BERT to verify our results. BERT (Kaliyar et al., 2021) and Bio-BERT (Gundapu & Mamidi, 2021) are good text classification algorithms. BERT is suitable for text classification, but it does not cover medical words. Bio-BERT corpora consist of text from PubMed abstracts (PubMed) and PubMed Central full-text articles (PMC). The text corpora of Bio-BERT consist of the wiki, Books, PubMed, and PMC (Lee et al., 2020).

The plot in figure 2 is generated after comparing the Tagging feeds. The precision obtained by all three is almost the same, but Recall, Accuracy, and F1-scores show that all three algorithms behaved differently. BERT is expected not to perform well on medical text as it does not have any medical dictionary in its corpus. Bio-BERT performed very well, but since the corpus is not updated dynamically, it misses out on a few new medical words added. Like Omicron, a variant of Covid-19 came up a few days back. This word would not be there in the pre-trained model. Our Text corpus is built manually and updated with the latest words when we run the code. TA-WHI also shows a marginal improvement of 1% over Bio-BERT. F1-score of TA-WHI and Bio-Bert show that they have more balanced Precision and Recall in comparison to BERT.

The focus of the current study is on Covid-19, heart attack, and Sugar.

The verification of the marked class is done using a WebCrawler we designed to check the contents on google. Keywords with the assigned class are searched; we mark them as TP, TN, FP, and FN, depending on the results collected. The accuracy in correctly recognizing the class related to COVID-19 is shown in Table 1. The average accuracy is 98.19 percent.



Figure 2. Performance of BERT, Bio-BERT, and TA-WHI

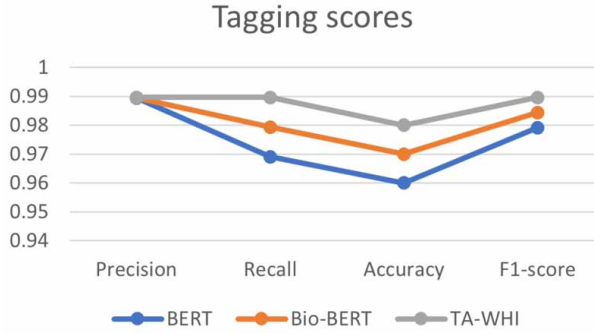


Table 1. Accuracy in correctly tagging COVID-19 feeds

Class	Precision	Recall	Accuracy
Sufficient	0.942256	1	98.37%
Misleading	0.93047	1	97.64%
Fabricated	0.9413	1	98.12%
Meaningful	0.9326	1	97.84%
Advertisement	0.9513	1	98.98%

Table 2 shows the accuracy of incorrectly identifying the Heart-Attack feeds’ classes. The average accuracy is 97.62%.

Table 3 computes the accuracy achieved by incorrectly identifying the class associated with Sugar. The average accuracy is 98.124%.

After the text was ranked and its class marked, CNN (Saigal & Khanna, 2020), LSTM (Yang et al., 2018), and CatBoost (Nasir et al., 2021) algorithms were applied to get the most robust predictive model. For pre-processing and tagging, we used Bio-BERT and TA-WHI with each algorithm. The results are reproduced in Table 4. LSTM is a deep learning algorithm and was expected to give better results. LSTM requires a larger dataset; our test results are based on 100 feeds collected over 120 days. An accuracy of 95% by LSTM is very close to the results with CatBoost. It is a lightweight boosting algorithm and can work well on smaller datasets. TA-WHI, with all the algorithms, performed better. An accuracy score of 98% with CatBoost is very encouraging.

Table 2. Accuracy in tagging of heart-attack feeds

Class	Precision	Recall	Accuracy
Sufficient	0.9234	1	96.76%
Misleading	0.9386	1	97.96%
Fabricated	0.9402	1	98.09%
Meaningful	0.9313	1	97.14%
Advertisement	0.9421	1	98.16%

**Table 3. Accuracy in tagging of sugar feeds**

Class	Precision	Recall	Accuracy
Sufficient	0.9342	1	97.45%
Misleading	0.9413	1	98.19%
Fabricated	0.9453	1	98.65%
Meaningful	0.9431	1	98.35%
Advertisement	0.9398	1	97.98%

**Table 4. Evaluation of different algorithms for classification**

Algorithm	Precision	Recall	Accuracy	F1-Score
CNN (+Bio-BERT)	0.945055	0.966292	0.92	0.972678
CNN (+TA-WHI)	0.978022	0.967391	0.95	0.972678
LSTM (+Bio-BERT)	0.988889	0.956989	0.95	0.972678
LSTM (+TA-WHI)	0.989247	0.978723	0.969697	0.983957
CatBoost (+Bio-BERT)	0.979381	0.979381	0.96	0.979381
CatBoost (+TA-WHI)	0.989796	0.989796	0.98	0.989796

## CONCLUSION

There is a Fake-BERT model also, but we preferred using Bio-BERT as it is more related to health information. Fake-BERT is good for classifying the general fake news, rumors, and misleading information, which is beyond the scope of our study. We used the pre-trained model of the classification algorithm and hence are not reproduced in the paper. The focus of TA-WHI is to correctly identify the class attached to the social media feeds. Different Machine learning algorithms are run over the refined dataset to create a predictive model. CNN can classify and predict the class with an accuracy of 95%, LSTM showed an accuracy of 96.86%, and CatBoost 98%. We can conclude from the results obtained that TA-WHI can accurately identify feeds as sufficient, fabricated, meaningful, advertisement, and misleading. Once TA-WHI is fully trained and tested, there will be no need to check the class of the feed using a third-party tool. The proposed model would predict the class associated with the social media feed. TA-WHI is not limited to health information. It can be used to correctly identify the emotion associated with any feed on the social media platform.

TA-WHI heavily depends on pre-processing to correctly identify the class of the feed. We have ignored emojis in the current research, but they can play an important role. In the future, we will convert the emojis to text. The classification algorithms work on vectors; we have used word2vec to convert the text into an equivalent number. Vocabulary expands every day, and there is always a chance of getting an inaccurate vector of the word. We would convert the words into vectors in the future using our algorithm. To improve the accuracy further, the authors would use a deep learning algorithm and create a dynamic benchmark dataset of Fake Health information on the web.

## AUTHOR'S CONTRIBUTION

The unique features of the proposed work are:

1. We designed a dictionary – MeDF, to retain relevant words from the feed.
2. We formulated polarity for the fake news related to health.

3. We designed a ranking method to accurately rank the feed based on the polarity.
4. We designed a classification formulation to mark the feed as “fabricated, misleading, sufficient, Irrelevant, and advertisement.” Classification is based on the ranking method.
5. The proposed work can be used to classify any dataset with slight modifications.

### **CONFLICT OF INTEREST**

There are no personal ties among the authors that may have seemed to affect the work presented in this study. Hence there is no conflict of interest between the authors.

## REFERENCES

- Agley, J., & Xiao, Y. (2021). Misinformation about COVID-19: Evidence for differential latent profiles and a strong association with trust in science. *BMC Public Health*, 21(1), 1–12. doi:10.1186/s12889-020-10103-x PMID:33413219
- Ahmed, A. A. A., Aljabouh, A., Donepudi, P. K., & Choi, M. S. (2021). Detecting fake news using machine learning: A systematic literature review. ArXiv Preprint ArXiv:2102.04458
- Alharbi, J. R., & Alhalabi, W. S. (2020). Hybrid approach for sentiment analysis of twitter posts using a dictionary-based approach and fuzzy logic methods: Study case on cloud service providers. [IJSWIS]. *International Journal on Semantic Web and Information Systems*, 16(1), 116–145. doi:10.4018/IJSWIS.2020010106
- Alrubaian, M., Al-Qurishi, M., Alamri, A., Al-Rakhami, M., Hassan, M. M., & Fortino, G. (2018). Credibility in online social networks: A survey. *IEEE Access: Practical Innovations, Open Solutions*, 7, 2828–2855. doi:10.1109/ACCESS.2018.2886314
- Andrea, K. (2022). *Text Analytics & NLP in Healthcare: Applications & Use Cases*. Lexalytics. <https://www.lexalytics.com/blog/text-analytics-nlp-healthcare-applications/>
- Atanasova, P., Nakov, P., Márquez, L., Barrón-Cedeño, A., Karadzhev, G., Mihaylova, T., Mohtarami, M., & Glass, J. (2019). Automatic fact-checking using context and discourse information. [JDIQ]. *Journal of Data and Information Quality*, 11(3), 1–27. doi:10.1145/3297722
- Bhattarai, B., Granmo, O.-C., & Jiao, L. (2021). Explainable tsetlin machine framework for fake news detection with credibility score assessment. ArXiv Preprint ArXiv:2105.09114.
- Burel, G., Farrell, T., & Alani, H. (2021). Demographics and topics impact on the co-spread of COVID-19 misinformation and fact-checks on Twitter. *Information Processing & Management*, 58(6), 102732. doi:10.1016/j.ipm.2021.102732 PMID:34511703
- Chawla, Y., Radziwon, A., Scaringella, L., Carlson, E. L., Greco, M., Silveira, P. D., de Aguiar, E. P., Shen, Q., Will, M., & Kowalska-Pyzalska, A. (2021). Predictors and outcomes of individual knowledge on early-stage pandemic: Social media, information credibility, public opinion, and behaviour in a large-scale global study. *Information Processing & Management*, 58(6), 102720. doi:10.1016/j.ipm.2021.102720
- Chen, E. S. (2016). Data, Text, and Web Mining in Healthcare. In L. Liu & M. T. Özsu (Eds.), *Encyclopedia of Database Systems* (pp. 1–5). Springer New York. doi:10.1007/978-1-4899-7993-3\_94-2
- Ciampaglia, S., Shiralkar, P., Rocha, L. M., Bollen, J., Menczer, F., & Flammini, A. (2015). Ciampaglia GL, Shiralkar P, Rocha LM, Bollen J, Menczer F, Flammini A. *Computational Fact Checking from Knowledge Networks*. *PLoS One*, 10(6), 1–13. doi:10.1371/journal.pone.0128193 PMID:26083336
- Ehrenstein, V., Kharrazi, H., & Lehmann, H. (2022). *Obtaining Data From Electronic Health Records - Tools and Technologies for Registry Interoperability, Registries for Evaluating Patient Outcomes: A User's Guide, 3rd Edition, Addendum 2 - NCBI Bookshelf*. <https://www.ncbi.nlm.nih.gov/books/NBK551878/>
- FactCheck. (2022). *A Project of The Annenberg Public Policy Center*. FactCheck. <https://www.factcheck.org/>
- FACTLY. (2022). *Making Public Data Meaningful*. FACTLY. <https://factly.in/>
- Fernández-Celemín, L., & Jung, A. (2006). What should be the role of the media in nutrition communication? *British Journal of Nutrition*, 96(S1), S86–S88. doi:10.1079/BJN20061707 PMID:16923259
- Full Fact. (2022). *Full Fact*. Full fact. <https://fullfact.org/>
- García-Peñalvo, F. J., Peraković, D., Mishra, A., & Gupta, B. B. (2021). *A Survey on Data mining classification approaches*. CEUR.
- Gundapu, S., & Mamidi, R. (2021). Transformer based automatic COVID-19 fake news detection system.
- Gunti, P., Gupta, B. B., & Benkhelifa, E. (2022). Data mining approaches for sentiment analysis in online social networks (OSNs). In *Data mining approaches for big data and sentiment analysis in social media* (pp. 116–141). IGI Global. doi:10.4018/978-1-7998-8413-2.ch005

- Harsin, J. (2018). Post-truth and critical communication studies. In Oxford research encyclopedia of communication. doi:10.1093/acrefore/9780190228613.013.757
- Kaliyar, R. K., Goswami, A., & Narang, P. (2021). FakeBERT: Fake news detection in social media with a BERT-based deep learning approach. *Multimedia Tools and Applications*, 80(8), 11765–11788. doi:10.1007/s11042-020-10183-2 PMID:33432264
- Kim, M., Gupta, B. B., & Rho, S. (2018). Crowdsourcing based scientific issue tracking with topic analysis. *Applied Soft Computing*, 66, 506–511. doi:10.1016/j.asoc.2017.09.028
- Kumari, R., Ashok, N., Ghosal, T., & Ekbal, A. (2021). Misinformation detection using multitask learning with mutual learning for novelty detection and emotion recognition. *Information Processing & Management*, 58(5), 102631. doi:10.1016/j.ipm.2021.102631
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., & Kang, J. (2020). Biobert: Pre-trained biomedical language representation model for biomedical text mining.
- Link, E. (2021). Information avoidance during health crises: Predictors of avoiding information about the COVID-19 pandemic among german news consumers. *Information Processing & Management*, 58(6), 102714. doi:10.1016/j.ipm.2021.102714 PMID:34539039
- Liu, L., & Tamer Özsu, M. (2018). Encyclopedia of Database Systems. Encyclopedia of Database Systems. doi:10.1007/978-1-4614-8265-9
- Mohammed, S. S., Menaouer, B., Zohra, A. F. F., & Nada, M. (2022). Sentiment Analysis of COVID-19 Tweets Using Adaptive Neuro-Fuzzy Inference System Models. [IJSSCI]. *International Journal of Software Science and Computational Intelligence*, 14(1), 1–20. doi:10.4018/IJSSCI.300361
- Murayama, T. (2021). Dataset of fake news detection and fact verification. *Survey (London, England)*.
- Nasir, J. A., Khan, O. S., & Varlamis, I. (2021). Fake news detection: A hybrid CNN-RNN based deep learning approach. *International Journal of Information Management Data Insights*, 1(1), 100007. doi:10.1016/j.jjimei.2020.100007
- PolitiFact. (2022). *PolitiFact*. PolitiFact. <https://www.politifact.com/>
- Potthast, M., Kiesel, J., Reinartz, K., Bevendorff, J., & Stein, B. (2017). A stylometric inquiry into hyperpartisan and fake news.
- Przybyła, P., & Soto, A. J. (2021). When classification accuracy is not enough: Explaining news credibility assessment. *Information Processing & Management*, 58(5), 102653. doi:10.1016/j.ipm.2021.102653
- Ranganath, S., Wang, S., Hu, X., Tang, J., & Liu, H. (2017). Facilitating time critical information seeking in social media. *IEEE Transactions on Knowledge and Data Engineering*, 29(10), 2197–2209. doi:10.1109/TKDE.2017.2701375
- Reis, J. C. S., Correia, A., Murai, F., Veloso, A., & Benevenuto, F. (2019). Explainable machine learning for fake news detection. *Proceedings of the 10th ACM Conference on Web Science*, 17–26. doi:10.1145/3292522.3326027
- Saigal, P., & Khanna, V. (2020). Multi-category news classification using Support Vector Machine based classifiers. *SN Applied Sciences*, 2(3), 1–12. doi:10.1007/s42452-020-2266-6
- Samadi, M., Mousavian, M., & Momtazi, S. (2021). Deep contextualized text representation and learning for fake news detection. *Information Processing & Management*, 58(6), 102723. doi:10.1016/j.ipm.2021.102723
- Sharma, Y., Bhargava, R., & Tadikonda, B. V. (2021). Named Entity Recognition for Code Mixed Social Media Sentences. [IJSSCI]. *International Journal of Software Science and Computational Intelligence*, 13(2), 23–36. doi:10.4018/IJSSCI.2021040102
- Singh, S. K., & Sachan, M. K. (2021). Classification of code-mixed bilingual phonetic text using sentiment analysis. [IJSWIS]. *International Journal on Semantic Web and Information Systems*, 17(2), 59–78. doi:10.4018/IJSWIS.2021040104
- Tacchini, E., Ballarin, G., della Vedova, M. L., Moret, S., & de Alfaro, L. (2017). Some like it hoax: Automated fake news detection in social networks.

The Washington Post. (2022). *Fact checker - The Washington Post*. The Washington Post. <https://www.washingtonpost.com/news/fact-checker/>

Thorne, J., Vlachos, A., Cocarascu, O., Christodoulopoulos, C., & Mittal, A. (2018). *The fact extraction and verification (fever) shared task*. ACL Anthology. doi:10.18653/v1/W18-5501

Veluru, S., Rahulamathavan, Y., Gupta, B. B., & Rajarajan, M. (2015). Privacy preserving text analytics: research challenges and strategies in name analysis. *Standards and Standardization: Concepts, Methodologies, Tools, and Applications*, 1415–1435.

Wang, X., Li, Y., Li, J., Liu, Y., & Qiu, C. (2021). A rumor reversal model of online health information during the Covid-19 epidemic. *Information Processing & Management*, 58(6), 102731. doi:10.1016/j.ipm.2021.102731 PMID:34539040

Xiong, X., Yang, B., & Kang, Z. (2018). A gradient tree boosting based approach to rumor detecting on sina weibo. arXiv preprint arXiv:1806.06326

Yang, Y., Zheng, L., Zhang, J., Cui, Q., Li, Z., & Yu, P. S. (2018). TI-CNN: Convolutional neural networks for fake news detection. arXiv preprint arXiv:1806.00749

Yao, Z., Tang, P., Fan, J., & Luan, J. (2021). Influence of online social support on the Public's belief in overcoming COVID-19. *Information Processing & Management*, 58(4), 102583. doi:10.1016/j.ipm.2021.102583 PMID:33746338

Zhang, Z., Sun, R., Wang, X., & Zhao, C. (2017). A situational analytic method for user behavior pattern in multimedia social networks. *IEEE Transactions on Big Data*, 5(4), 520–528. doi:10.1109/TBDATA.2017.2657623

Zhou, C., Li, K., & Lu, Y. (2021). Linguistic characteristics and the dissemination of misinformation in social media: The moderating effect of information richness. *Information Processing & Management*, 58(6), 102679. doi:10.1016/j.ipm.2021.102679

Zhou, C., Xiu, H., Wang, Y., & Yu, X. (2021). Characterizing the dissemination of misinformation on social media in health emergencies: An empirical study based on COVID-19. *Information Processing & Management*, 58(4), 102554. doi:10.1016/j.ipm.2021.102554 PMID:36570740

Zubiaga, A., Liakata, M., Procter, R., Wong Sak Hoi, G., & Tolmie, P. (2016). Analysing How People Orient to and Spread Rumours in Social Media by Looking at Conversational Threads. *PLoS One*, 11(3), e0150989. doi:10.1371/journal.pone.0150989 PMID:26943909

Zwolenski, M., & Weatherill, L. (2014). The digital universe: Rich data and the increasing value of the Internet of things. *Journal of Telecommunications and the Digital Economy*, 2(3), 41–47. doi:10.7790/ajtde.v2n3.47

*Piyush Bagla is pursuing his Ph.D. at Dr. B.R. Ambedkar National Institute of Technology, Jalandhar, Punjab, in the Dept. of Computer Science and Engineering. His research area includes Machine Learning (ML) and Natural Language Processing (NLP). Prior to joining NITJ, he worked as an Assistant Professor at Graphic Era Hill University (GEHU), Dehradun, Uttarakhand. He pursued his Master of Technology at Graphic Era deemed to be University, Dehradun, Uttarakhand, and his Bachelor of Technology (B.Tech) from Uttarakhand Technical University (UTU). He has good research publication experience as well.*

*Kuldeep Kumar is currently working as an Assistant Professor in the Department of Computer Engineering at National Institute of Technology Kurukshetra. He pursued his Ph.D. from the National University of Singapore (NUS), Singapore. His research interests are Software Engineering, Machine Learning, and Data Analytics. He has published more than 40 papers in peer-reviewed journals and international conferences along with several book chapters.*