

# Fake News Detection of South African COVID-19 Related Tweets using Machine Learning

Yaseen Khan and Surendra Thakur  
*Department of Information Technology  
 Durban University of Technology  
 Durban, South Africa  
 khanyas786@gmail.com, thakur@dut.ac.za*

**Abstract**— Social Media has grown in popularity in recent years comprising of billions of users who in turn exchange and communicate content at a volume and rate impractical to examine manually. Fake News are now being used on these platforms to manipulate and affect societies across the world as was the case in the 2016 United States of America (USA) elections and of recent during the 2019 coronavirus disease (COVID-19) pandemic. South Africa is not immune to the spread of Fake News, particularly, through Social Media platforms such as Twitter, Facebook and TikTok. It is, therefore, important to detect the presence of Fake News computationally in order assist the mitigation of its spread and prevent perceivable negative effects. This study addresses the issue by developing a Machine Learning (ML) model to analyze large amounts of data associated with Social Media. Curated annotated datasets from CONSTRAINT AAAI 2021; COVID-19 Rumour, FNIR and Zenodo's COVID-19 datasets; Google and Polifact Fact Checked websites; were utilized to develop the ML model. Specifically, the model was trained on 36254 data points and applied on a South African related COVID-19 Twitter dataset collected for cursory analysis. In total, 27 ML models were experimented with and the collected South African COVID-19 related Twitter dataset comprised of 976087 tweets from 8 November 2020 until 19 July 2021. The results detected 329107 tweets as being 'Fake' based on the LightGBM Classifier which was chosen as the most feasible model in terms of speed and a balanced accuracy score of 0.82. The proposed model is unique as it is trained on a larger combined dataset and supplements existing efforts to combat misinformation, disinformation and malinformation spread on Twitter.

**Keywords**—Machine Learning, fake news, covid-19, NLP

## I. INTRODUCTION

Social Media has rapidly progressed in recent times such that information is now generated and transmitted at unprecedented high volumes and speed [1]. This has resulted in the production of unstructured Big Data which arrives and introduces certain challenges in each of its widely recognized constituents commonly referred to as 5 V's namely, Volume, Velocity, Variety, Value and Veracity [1]-[4]. Veracity, which analyses the truthfulness or factualness of data, is key not only for creating quality datasets but also to monitor trends, patterns and characteristics of influential role-players who generate incorrect information. The unintentional or purposeful generation of incorrect information on Social Media may be detrimental to naïve as well as experienced users who may be manipulated by being misinformed [5]. This study focuses on the veracity component of the 5 V's in terms of Fake News.

Fake News is a form of news related information that is incorrect or misleading and has been known to be used during election campaigns, pandemics and on a variety of subjects [6-9]. Fake News may be considered a sub-category of misinformation, with the latter encompassing non-news related information. Recently, Facebook, Twitter and other popular Social Media platforms have been criticised for their perceived lack of effort in combating against content that contain Fake News as well as incitement [10]-[12]. Consequently, these Social Media platforms are in the process of reviewing policies [13] and funding research [14] on Fake News, particularly misinformation, malinformation and disinformation. Malinformation is fake information to make a scheme sound more believable while disinformation is intended to cause harm or damage by using incorrect information [15].

Social Media data comprises of unstructured data in the form of audio, text, images, and videos. This unstructured data is complex to efficiently analyse and require specialized techniques which can be found in a multi-disciplinary field known as Data Science. Data Science is a main driver of growing advances in Artificial Intelligence and provides intelligent insights by computationally tackling issues surrounding Big Data and Social Media [16]-[19].

This study leveraged common machine learning techniques and developed a new fake news detection model built upon previous datasets and applied it on a South African Twitter related dataset for cursory analysis. The next section discusses related literature followed by the methodology, results and findings, conclusion, limitations and future work.

## II. RELATED LITERATURE

### A. COVID-19

Coronavirus disease 2019 (COVID-19) is a disease caused by a new type of coronavirus known as SARS-CoV-2 which was identified following a December 2019 outbreak in China [20]-[22]. The World Health Organisation (WHO) had officially declared COVID-19 a pandemic on 11 March 2020 [23] which had a global impact on numerous fields such as health, economy and politics. Given the high profile nature of COVID-19, it is worth exploring its social dialogue amongst communities, particularly in South Africa, with a particular focus on Fake News since its negative effects can range from minimal to fatal. Twitter, a popular Social Media platform, is an appropriate source to leverage for the exploration of Fake News discussions surrounding COVID-19 since it has already

been used in many studies such as sentiment analysis [24], [25] hate speech [26], [27] and rumour detection [28], [29].

### B. Social Media and Fake News

Social Media are Web 2.0 Internet-based applications designed for user-generated content which facilitates the development of online social networks [30]. Social Media's popularity and the growth of its users are significant, with 2.86 billion users in 2017 and 4.2 billion users in 2021, and with an average usage of 145 minutes a day per user [31]. These statistics reflect the growing concern that Social Media may be used as a medium to widely manipulate and spread propaganda like the case of the United States of America (USA) elections in 2016. It was researched that Fake News were widely deployed on Social Media to sway votes and manipulate the elections [6]. Fake News on Social Media also affected South Africa with the government creating websites to warn against fake COVID-19 related news on Twitter, such as the loan free debt relief on Facebook and the non-social distance school gathering crowds in the North West [32].

Fake News has often been categorised into Misinformation, Disinformation and Malinformation [15]. The Council of Europe [33] has defined these categorisations as follows: Disinformation - information that is false and deliberately created to harm a person, social group, organisation or country; Misinformation - information that is false, but not created with the intention of causing harm; Malinformation - information that is based in reality, but is used to inflict harm on a person, organisation or country. Fake News has been known to be leveraged by certain users on Social Media platforms in order to deceive, manipulate and spread false content at a large scale. Ahmed et al. [34] have highlighted how Social Media platforms such as Twitter can be utilised for monitoring public views and opinions related to disease outbreaks such as swine flu, Ebola, and the Zika virus. This can potentially support health authorities to discern potential misinformation, understand public concerns, and collect unanswered questions.

### C. Machine Learning

Machine learning (ML) is a technique considered important to solving many problems, which utilises appropriate algorithms and has broad applications [35]. ML techniques are frequently classified into either Supervised Learning or Unsupervised Learning. There is also a third technique referred to as Semi-supervised Learning. In Supervised Learning a computer is programmed to compute a function  $y = f(x)$ , where  $x$  is the input value, which is a set of examples supplied by humans also called the training data and  $y$  is the corresponding expected output value [24]. The function  $(x)$  thereafter establishes a pattern to produce output values  $y$  based on unseen input values  $x$  and this is referred to as computer learning. This technique is therefore classified as 'supervised' since the output  $y$  is guided by the value  $x$ , which is determined by humans thus creating a supervision scenario for the computer. Unsupervised Learning works with samples that are unlabelled [35]. Frequent unsupervised algorithms include clustering models such as k-means, k-medoids and hierarchical as well as hidden Markov models and certain unsupervised neural network models such as self-organising maps [36].

ML has been widely used as a classification approach in various Social Media studies to accomplish many tasks such as Sentiment Analysis, Hate Speech Detection, Emotional analysis and Fake News detection [37]-[39]. Deep Learning is another approach sometimes used but has some cost disadvantages. It usually requires more computational power associated with more expensive hardware compared to conventional ML models.

### D. Fake News Detectors

There have been several attempts to detect Fake News content with large volumes of literature incorporating several Data Science approaches. Bondielli and Marcelloni [28] surveyed Fake News and rumour detection techniques and provided varying techniques adopted by the researchers which includes ML, Deep Learning and other approaches. In terms of specifically detecting COVID-19 fake news various systems and models have been developed such as Cross-SEAN [40], CoVerifi [41] and Checkovid [42]. These models incorporate transformer, neural network, natural language processing (NLP) and ML techniques in order to computationally analyse and classify information.

## III. METHODOLOGY

### A. Data collection and processing

#### 1) Training dataset

The dataset used to train, test and validate the ML models was constructed using data from CONSTRAINT AAAI-2021, COVID-19 Rumour dataset, Zenodo, COVID-19 Fake News Infodemic Research Dataset (FNIR), Google Fact Check and Politifact websites [43]-[48]. An example can be seen in Table 1. The data was processed to remove duplicates, URLs, stopwords, emojis, and null entries. Furthermore, the classification labels were computationally transformed into 1's and 0's. Labels that contained 'FALSE', 'misleading', 'F', '0' and 'fake' were transformed into 1. Labels that contained 'TRUE', 'T', 'U', '1' and 'real' were transformed into 0. The reason for the transformation was to ensure consistency amongst labels. For this research, 1 represents potential fake news and 0 is not fake news. A total of 36254 data points remained consisting of text and label columns<sup>1</sup> with 24173 and 12082 representing 'Fake' and 'Not Fake' news respectively. The reason to combine training datasets from several sources was to create a more data rich training set which will improve ML performance.

Table 1 Example of training dataset

Text	Label
The NZ COVID Tracker app will remain important and useful at Alert level 1. People are encouraged to download the app or a similar one or keep a record of where you're going to be.	Real
BREAKING NEWS# The president Cyril Ramaphosa has asked all foreign nations to depart south Africa before 21 june 2020 due to increasing cases of COVID 19	Fake
We are delighted that 78 high- and upper-middle countries and economies have now confirmed they will participate in the COVAX Facility and the number is growing. I urge those who have not yet joined to do so by the 18th of September-@DrTedros #COVID19	Real

<sup>1</sup><https://github.com/khanYas786/DUT/blob/master/Text-label-xl-cleaned.csv>

## 2) South African related COVID-19 Twitter dataset

In order to obtain a SA COVID-19 related Twitter dataset, data was collected from Twitter using a snscreape<sup>2</sup>. Snscreape is a social networking services (SNS) scraper and collects information such as user profiles, hashtags and tweets. The data was scraped using case-insensitive keywords, Covid19sa, Covid-19sa, Covid19-sa, coronavirussa, sacoronavirus, covidinsa, covid\_19sa, lockdownsa and covid19insa which were guided by the top tweeting trends [49]. The timeline is from 08 November 2020 until 19 July 2021 and has a total of 976087 tweets. The reason for these keywords was to focus on South African related content rather than global content which would have been unfeasible to process due to the many tweets associated with the COVID-19 pandemic. The timeline of collection was based on the researchers' initiation of the study. The metadata includes 29 columns such as the ID, Content, User and Date fields.

### B. Model training and evaluation

Each ML process split the modeling data into 75% for training and 25% for testing. This is a common ratio amongst research and data science practitioners. Due the imbalanced nature of the modeling dataset, the 'Balanced Accuracy' performance metric [50] was prioritized as the decision making parameter. In total, 27 models were experimented with using Lazy Predict<sup>3</sup>. Lazy Predict is a Python package that builds several basic ML models without any parameter tuning and provides information on which models performance best.

## IV. FINDINGS AND RESULTS

The performance metrics of the experimented 27 ML models utilized can be seen in Table 2. It is sorted by highest to lowest based on the Balanced Accuracy metric with the best performing model listed at the top. The Extra-Trees Classifier model outperforms the rest but takes longer to process than the Light Gradient Boosting Machine (LightGBM) Classifier. LightGBM is a gradient boosting framework that uses tree-based learning algorithms [51]. Due to the SA COVID-19 dataset being large it was decided to choose the LightGBM Classifier as the desired model to parse through the data.

Table 2 Results of ML models evaluated

Model	Accuracy	Balanced Accuracy	ROC AUC	F1 Score	Time Taken
ExtraTreesClassifier	0.869	0.823	0.823	0.864	88.635
LightGBMClassifier	0.860	0.820	0.820	0.856	7.375
SVC	0.862	0.816	0.816	0.857	1881.305
RandomForestClassifier	0.862	0.815	0.815	0.856	46.048
XGBClassifier	0.860	0.815	0.815	0.854	18.694
LogisticRegression	0.836	0.801	0.801	0.833	2.453
BernoulliNB	0.832	0.797	0.797	0.829	2.372
LinearDiscriminantAnalysis	0.838	0.797	0.797	0.834	9.030
BaggingClassifier	0.844	0.795	0.795	0.837	142.967
RidgeClassifierCV	0.838	0.795	0.795	0.833	7.438

<sup>2</sup><https://github.com/JustAnotherArchivist/snscreape>

<sup>3</sup><https://github.com/shankarpandala/lazypredict>

RidgeClassifier	0.838	0.795	0.795	0.833	2.237
NearestCentroid	0.819	0.793	0.793	0.818	1.934
DecisionTreeClassifier	0.812	0.782	0.782	0.811	29.889
QuadraticDiscriminantAnalysis	0.789	0.781	0.781	0.792	7.877
CalibratedClassifierCV	0.833	0.781	0.781	0.825	342.846
SGDClassifier	0.810	0.780	0.780	0.808	1651.793
NuSVC	0.840	0.775	0.775	0.828	
LinearSVC	0.803	0.774	0.774	0.802	88.949
AdaBoostClassifier	0.825	0.774	0.774	0.818	25.992
ExtraTreeClassifier	0.797	0.760	0.760	0.794	2.744
Perceptron	0.776	0.752	0.752	0.777	2.180
GaussianNB	0.738	0.744	0.744	0.744	2.386
PassiveAggressiveClassifier	0.764	0.737	0.737	0.764	2.921
KNeighborsClassifier	0.721	0.598	0.598	0.661	446.224
LabelSpreading	0.667	0.510	0.510	0.542	213.344
LabelPropagation	0.667	0.510	0.510	0.542	106.224
DummyClassifier	0.558	0.499	0.499	0.554	1.861

Figure 1 shows the confusion matrix of the LightGBM Classifier.

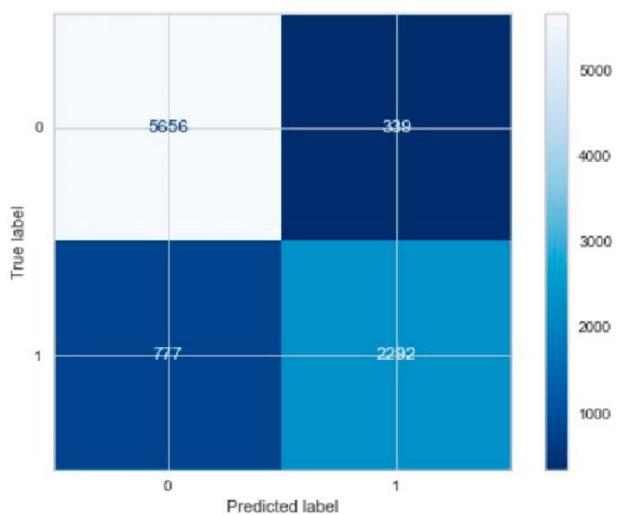


Figure 1 Confusion matrix of LightGBM Classifier

The parameters for the LightGBM model were as follows: Boosting type = 'Gradient-boosted decision tree'; Class weight = None; Colsample bytree = 1.0; Importance type='split'; Learning rate = 0.09; Max depth = -5; Min child samples=20; Min child weight=0.001; Min split gain=0.0; n estimators = 440; n jobs = -1; num leaves = 31; objective = None; random state = 42; reg alpha = 0.0; reg lambda = 0.0.; silent=True; subsample = 1.0; subsample for bin=200000; subsample freq=0.

The South African COVID-19 related dataset<sup>4</sup> was parsed through the LightGBM model and detected 329107 tweets as

containing fake news and 646981 tweets as not being fake. Figure 2 depicts a longitudinal trend of 'Fake' and 'Not Fake' news content detected from the collected South African COVID-19 dataset.

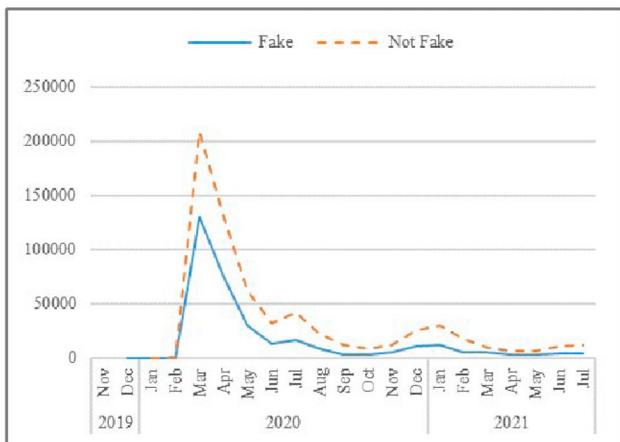


Figure 2 'Fake' and 'Not Fake' news frequency trends of the South African COVID-19 Twitter Dataset

## V. DISCUSSION AND CONCLUSION

The study provided a new Fake News Detection model for COVID-19 by combining data from recent research and supplementary material. A satisfactory model with a 'Balanced Accuracy' of 82% was developed which provided the means to detect Fake News on a South African related Twitter dataset for cursory analysis. This model was selected due to it being faster than other related classifiers and can computationally analyze large amounts of data in a short period of time, consequently, aiding efforts in combating fake news on social media.

## VI. LIMITATIONS OF THE STUDY

Extracting data from Twitter may vary depending on a time point. Furthermore Twitter content may be removed, deleted or made private by users or the platform itself. The methods used does not reflect all ML algorithms and cannot be applied on audio, images or videos. In addition, due to hardware and costing constraints models involving transformer or neural network techniques could not be explored satisfactorily on the large dataset.

## VII. FUTURE WORK

Deep learning and NLP approaches, although often hardware intensive, should be explored which may provide more accurate fake news detection models. In addition, since social media contain data such as audio, images and videos, a multi-modal design that can detect fake news across all formats need to be investigated.

## REFERENCES

- [1] M. Chen, S. Mao and Y. Liu, "Big Data: A Survey", *Mobile Networks and Applications*, vol. 19, no. 2, pp. 171-209, 2014. Available: 10.1007/s11036-013-0489-0.
- [2] M. Khan, M. Uddin and N. Gupta, "Seven V's of Big Data understanding Big Data to extract value", *Proceedings of the 2014 Zone I Conference of the American Society for Engineering Education*, pp. 1-5, 2014. Available: 10.1109/aseezone1.2014.6820689 [Accessed 31 March 2022].
- [3] G. Bello-Orgaz, J. Jung and D. Camacho, "Social big data: Recent achievements and new challenges", *Information Fusion*, vol. 28, pp. 45-59, 2016. Available: 10.1016/j.inffus.2015.08.005.
- [4] G. Manogaran, C. Thota, D. Lopez, V. Vijayakumar, K. Abbas and R. Sundarsekhar, "Big Data Knowledge System in Healthcare", *Studies in Big Data*, pp. 133-157, 2017. Available: 10.1007/978-3-319-49736-5\_7 [Accessed 31 March 2022].
- [5] M. Visentini, G. Pizzi and M. Pichierri, "Fake News, Real Problems for Brands: The Impact of Content Truthfulness and Source Credibility on consumers' Behavioral Intentions toward the Advertised Brands", *Journal of Interactive Marketing*, vol. 45, pp. 99-112, 2019. Available: 10.1016/j.intmar.2018.09.001.
- [6] H. Allcott and M. Gentzkow, "Social Media and Fake News in the 2016 Election", *Journal of Economic Perspectives*, vol. 31, no. 2, pp. 211-236, 2017. Available: 10.1257/jep.31.2.211.
- [7] A. Gelfert, "Fake News: A Definition", *Informal Logic*, vol. 38, no. 1, pp. 84-117, 2018. Available: 10.22329/il.v38i1.5068.
- [8] D. Lazer et al., "The science of fake news", *Science*, vol. 359, no. 6380, pp. 1094-1096, 2018. Available: 10.1126/science.aaq2998 [Accessed 31 March 2022].
- [9] K. Sharma, F. Qian, H. Jiang, N. Ruchansky, M. Zhang and Y. Liu, "Combating Fake News", *ACM Transactions on Intelligent Systems and Technology*, vol. 10, no. 3, pp. 1-42, 2019. Available: 10.1145/3305260.
- [10] Euronews, "Big Tech CEOs face questions in U.S. Congress over misinformation", *euronews*, 2022. [Online]. Available: <https://www.euronews.com/2021/03/25/facebook-twitter-and-google-face-questions-in-us-congress-over-misinformation>. [Accessed: 31-Mar-2022].
- [11] J. Wakefield, "Google, Facebook Twitter grilled in US on fake news", *BBC News*, 2022. [Online]. Available: <https://www.bbc.com/news/technology-56523378>. [Accessed: 31-Mar-2022].
- [12] BBC News, "Facebook, Twitter and Google face questions from US senators", *BBC News*, 2020 [Online]. Available: <https://www.bbc.com/news/technology-54721023> [Accessed 19 April 2021].
- [13] K. Paul, "Zuckerberg: Facebook will review policies after backlash over Trump posts", *The Guardian*, 2020. [Online]. Available: <https://www.theguardian.com/technology/2020/jun/05/mark-zuckerberg-facebook-trump-policies-review> [Accessed 19 April 2021].
- [14] Facebook, "Foundational Integrity Research: Misinformation and Polarization request for proposals - Facebook Research, Facebook 2020", [Online]. Available: <https://research.fb.com/programs/research-awards/proposals/foundational-integrity-research-misinformation-and-polarization-request-for-proposals/> [Accessed 19 April 2021].
- [15] M. Elhadad, K. Fun Li and F. Gebali, "Fake News Detection on Social Media: A Systematic Survey", *2019 IEEE Pacific Rim Conference on Communications, Computers and Signal Processing (PACRIM)*, 2019. Available: 10.1109/pacrim47961.2019.8895062 [Accessed 31 March 2022].
- [16] S. Asur and B. Huberman, "Predicting the Future with Social Media", *2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, 2010. Available: 10.1109/wi-iat.2010.63 [Accessed 31 March 2022].
- [17] F. Provost and T. Fawcett, "Data Science and its Relationship to Big Data and Data-Driven Decision Making", *Big Data*, vol. 1, no. 1, pp. 51-59, 2013. Available: 10.1089/big.2013.1508.
- [18] Oracle, "What is data Science?", *Oracle*, 2020. [Online]. Available: <https://www.oracle.com/data-science/what-is-data-science.html> [Accessed 18 November 2020].
- [19] M. Cinelli et al., "The COVID-19 social media infodemic", *Scientific Reports*, vol. 10, no. 1, 2020. Available: 10.1038/s41598-020-73510-5.
- [20] W. Guan et al., "Clinical Characteristics of Coronavirus Disease 2019 in China", *New England Journal of Medicine*, vol. 382, no. 18, pp. 1708-1720, 2020. Available: 10.1056/nejmoa2002032 [Accessed 31 March 2022].
- [21] Z. Wu and J. McGoogan, "Characteristics of and Important Lessons From the Coronavirus Disease 2019 (COVID-19) Outbreak in China", *JAMA*, vol. 323, no. 13, p. 1239, 2020. Available: 10.1001/jama.2020.2648.

- [22] Z. Zu et al., "Coronavirus Disease 2019 (COVID-19): A Perspective from China", *Radiology*, vol. 296, no. 2, pp. E15-E25, 2020. Available: 10.1148/radiol.2020200490.
- [23] D. Cucinotta and M. Vanelli, "WHO Declares COVID-19 a Pandemic", *Acta Biomed*, vol. 91, no. 1, pp. 157–160, Mar. 2020.
- [24] H. Saif, F. Ortega, M. Fernández and I. Cantador, "Sentiment Analysis in Social Streams", *Human–Computer Interaction Series*, pp. 119-140, 2016. Available: 10.1007/978-3-319-31413-6\_7 [Accessed 2 April 2022].
- [25] Y. Khan, S. Thakur, O. Obiyemi and E. Adetiba, "Exploring Links between Online Activism and Real-World Events: A Case Study of the #FeesMustFall", *Scientific Programming*, vol. 2022, pp. 1-10, 2022. Available: 10.1155/2022/1562592 [Accessed 2 April 2022].
- [26] Z. Waseem and D. Hovy, "Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter", 2022.
- [27] H. Watanabe, M. Bouazizi and T. Ohtsuki, "Hate Speech on Twitter: A Pragmatic Approach to Collect Hateful and Offensive Expressions and Perform Hate Speech Detection", *IEEE Access*, vol. 6, pp. 13825-13835, 2018. Available: 10.1109/access.2018.2806394 [Accessed 4 April 2022].
- [28] A. Bondielli and F. Marcelloni, "A survey on fake news and rumour detection techniques", *Information Sciences*, vol. 497, pp. 38-55, 2019. Available: 10.1016/j.ins.2019.05.035.
- [29] L. Tian, X. Zhang, Y. Wang and H. Liu, "Early Detection of Rumours on Twitter via Stance Transfer Learning", *Lecture Notes in Computer Science*, pp. 575-588, 2020. Available: 10.1007/978-3-030-45439-5\_38 [Accessed 4 April 2022].
- [30] J. Obar and S. Wildman, "Social media definition and the governance challenge: An introduction to the special issue", *Telecommunications Policy*, vol. 39, no. 9, pp. 745-750, 2015. Available: 10.1016/j.telpol.2015.07.014.
- [31] S. Kemp, "Digital 2021: the latest insights into the 'state of digital' - We Are Social UK", *We Are Social UK*, 2022. [Online]. Available: <https://wearesocial.com/uk/blog/2021/01/digital-2021-the-latest-insights-into-the-state-of-digital/>. [Accessed: 04- Apr- 2022].
- [32] South African Government, "Fake news - Coronavirus COVID-19 | South African Government", *Gov.za*, 2022. [Online]. Available: <https://www.gov.za/covid-19/resources/fake-news-coronavirus-covid-19>. [Accessed: 04- Apr- 2022].
- [33] C. Wardle and H. Derakhshan, "Information disorder: Toward an interdisciplinary framework for research and policy making", *Council of Europe*, 2017. [Online]. Available: <https://rm.coe.int/information-disorder-toward-an-interdisciplinary-framework-for-research/168076277c>. [Accessed: 04- Apr- 2022].
- [34] W. Ahmed, P. Bath, L. Sbaffi and G. Demartini, "Novel insights into views towards H1N1 during the 2009 Pandemic: a thematic analysis of Twitter data", *Health Information & Libraries Journal*, vol. 36, no. 1, pp. 60-72, 2019. Available: 10.1111/hir.12247.
- [35] C. ZHAI and S. MASSUNG, *Text data management and analysis*. [New York, NY] ; [San Rafael, California]: Morgan and Claypool, 2016.
- [36] A. Wade and G. Di MarzoSerugendo, "A Model of Extracting Patterns in Social Network Data Using Topic Modelling, Sentiment Analysis and Graph Databases", *Computer Science & Information Technology (CS & IT)*, 2017. Available: 10.5121/csit.2017.70608 [Accessed 4 April 2022].
- [37] P. Yang and Y. Chen, "A survey on sentiment analysis by using machine learning methods", *2017 IEEE 2nd Information Technology, Networking, Electronic and Automation Control Conference (ITNEC)*, 2017. Available: 10.1109/itnec.2017.8284920 [Accessed 4 April 2022].
- [38] N. Aulia and I. Budi, "Hate Speech Detection on Indonesian Long Text Documents Using Machine Learning Approach", *Proceedings of the 2019 5th International Conference on Computing and Artificial Intelligence - ICCAI '19*, 2019. Available: 10.1145/3330482.3330491 [Accessed 4 April 2022].
- [39] J. Zhang, Z. Yin, P. Chen and S. Nichele, "Emotion recognition using multi-modal data and machine learning techniques: A tutorial and review", *Information Fusion*, vol. 59, pp. 103-126, 2020. Available: 10.1016/j.inffus.2020.01.011.
- [40] W. Paka, R. Bansal, A. Kaushik, S. Sengupta and T. Chakraborty, "Cross-SEAN: A cross-stitch semi-supervised neural attention model for COVID-19 fake news detection", *Applied Soft Computing*, vol. 107, p. 107393, 2021. Available: 10.1016/j.asoc.2021.107393.
- [41] N. Kolluri and D. Murthy, "CoVerifi: A COVID-19 news verification system", *Online Social Networks and Media*, vol. 22, p. 100123, 2021. Available: 10.1016/j.osnem.2021.100123.
- [42] S. Dadgar and M. Ghatee, "Checkovid: A COVID-19 misinformation detection system on Twitter using network and content mining perspectives", *arXiv.org*, 2022. [Online]. Available: <https://doi.org/10.48550/arXiv.2107.09768>. [Accessed: 04- Apr- 2022].
- [43] P. Patwa et al., "Fighting an Infodemic: COVID-19 Fake News Dataset", *Combating Online Hostile Posts in Regional Languages during Emergency Situation*, pp. 21-29, 2021. Available: 10.1007/978-3-030-73696-5\_3 [Accessed 4 April 2022].
- [44] M. Cheng et al., "A COVID-19 Rumor Dataset", *Frontiers in Psychology*, vol. 12, 2021. Available: 10.3389/fpsyg.2021.644801 [Accessed 4 April 2022].
- [45] S. Banik, "COVID Fake News Dataset", *Zenodo*, 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.4282522>. [Accessed: 04- Apr- 2022].
- [46] J. Saenz, S. Kalathur Gopal and D. Shukla, "Covid-19 Fake News Infodemic Research Dataset (CoVID19-FNIR Dataset)", *IEEE DataPort*, 2021. [Online]. Available: <https://dx.doi.org/10.21227/b5bt-5244>. [Accessed: 04- Apr- 2022].
- [47] Google, "Fact Check Tools", *Toolbox.google.com*. [Online]. Available: <https://toolbox.google.com/factcheck/explorer>. [Accessed: 04- Apr- 2022].
- [48] Politifact, "PolitiFact | Coronavirus", *Politifact.com*. [Online]. Available: <https://www.politifact.com/coronavirus/>. [Accessed: 04- Apr- 2022].
- [49] Media Update, "Seven trending hashtags about COVID-19 on social media", *Media Update*, 2020. [Online]. Available: <https://www.mediaupdate.co.za/social/148423/seven-trending-hashtags-about-covid-19-on-social-media>. [Accessed: 04- Apr- 2022].
- [50] K. Brodersen, C. Ong, K. Stephan and J. Buhmann, "The Balanced Accuracy and Its Posterior Distribution", *2010 20th International Conference on Pattern Recognition*, 2010. Available: 10.1109/icpr.2010.764 [Accessed 4 April 2022].
- [51] GeeksforGeeks, "LightGBM (Light Gradient Boosting Machine) - GeeksforGeeks", *GeeksforGeeks*, 2021. [Online]. Available: <https://www.geeksforgeeks.org/lightgbm-light-gradient-boosting-machine/>. [Accessed: 04- Apr- 2022].