# Text Analysis of COVID-19 Tweets

Panagiotis C. Theocharopoulos[1](✉) , Anastasia Tsoukala[1],
Spiros V. Georgakopoulos[2], Sotiris K. Tasoulis[1], and Vassilis P. Plagianakos[1]

[1] Department of Computer Science and Biomedical Informatics,
University of Thessaly, Lamia, Greece
{ptheochar,antsoukala,stasoulis,vpp}@uth.gr
[2] Department of Mathematics, University of Thessaly, Lamia, Greece
spirosgeorg@uth.gr

**Abstract.** During the COVID-19 pandemic many countries were forced to implement lockdowns to prevent further spread of the SARS-CoV-2, prohibiting people from face-to-face social interactions. This unprecedented circumstance led to an increase in traffic on social media platforms, one of the most popular of which is Twitter, with a diverse spectrum of users from around the world. This quality, along with the ability to use its API for research purposes, makes it a valuable resource for data collection and analysis. In this paper we aim to present the sentiments towards the COVID-19 pandemic and vaccines as it was imprinted through the users' tweets when the events were actually still in motion. For our research, we gathered the related data from Twitter and characterized the gathered tweets in two classes, positive and negative; using the BERT model, with an accuracy of 99%. Finally, we performed various time series analyses based on people's sentiment with reference to the pandemic period of 2021, the four major vaccine's companies as well as on the vaccine's technology.

**Keywords:** Text classification · COVID-19 · Vaccines · BERT · Sentiment analysis · Twitter

## 1 Introduction

A major pivot in recent human history, specifically in 2020, was boldly highlighted by the COVID-19 global pandemic. Shortly thereafter the year of 2021 is unquestionably marked by the pharmaceutical attempt to shield the masses against the virus and its variants. The medical scientists seeking out the prevention of the disease were not the only researchers working to prevent a societal fall. The precariousness of the pandemic set in motion a big part of the research community to search for patterns of certainty within the data in order to level the global playing field [22]. Twitter attracted them, as it is a major source for tagged and specified data aggregation.

At the beginning of the vaccination period citizens were more agreeable to it, mainly for the purpose of achieving herd immunity against the highly infectious virus as soon as possible, and the realisation of the returning-to-normality

promise [8,15]. However, a significant minority of citizens in many countries expressed a general refusal to vaccination either due to mistrust to the local and global authorities or distrust to the effectiveness, safety and necessity of the COVID-19 vaccines [2].

## 1.1    COVID-19 Vaccines

Leading companies in the pharmaceutical industry attempted to administrate the plans of reconciliation with the help of financial and advertising support of the governments [22]. In some countries, mandates were actually applied with the result of further polarizing the relationships between the governments and the people [6,16]. A geopolitical division based on major vaccine companies also took place, as more vaccines were produced. However, not all of them were globally authorised or approved for emergency administration [3]. The four major vaccine companies that managed to obtain a standing public image, at least in the English speaking world, in the entire Europe and in some African, Oceanian and Asian countries are AstraZeneca, Johnson & Johnson, Moderna and Pfizer. Milken Institute's website about COVID-19 Vaccine Tracker showed on the 18th of December 2020 that 236 vaccines were in development, 38 were then in clinical testing and 7 had reached a regulatory decision [4]. In just a year of clinical trials and regulatory authorization protocols, the four leading vaccines have today managed to get global emergency authorization, with the exception of Cuba, Venezuela and some parts of Africa and Asia. We observe that Twitter and the four leading vaccines share a common factor; the countries that mostly use Twitter are also the ones that have either emergency use or full authorization of these vaccines [3].

## 1.2    Twitter

The Twitter website started as a patchwork infrastructure and organically progressed into an integrating one [14]. Twitter relies on open source tools and a simplistic design to maintain the micro-blogging of its users [14]. The company brought up the innovation of hashtags, when no other social media could filter its content so profoundly. Of course, this asset is its monetization tool for stakeholders and advertisers. In the "Global Impact Report" posted on their website on 2021 Twitter mentions that during the pandemic governments and vaccine companies were amongst various stakeholders contributing to decision-making and the advertising of vaccine campaigns[1]. Among actions taken from 2020 until March 2022 concerning the platform's and its stakeholders' initiative fighting misinformation during the pandemic, Twitter reported about 6,7 thousand account suspensions and 78 thousand content removals[2]

---

[1] https://about.twitter.com/content/dam/about-twitter/en/company/global-impact-2020.pdf.

[2] Twitter COVID-19 Misinformation Report.

### 1.3   Sentiment Analysis

The application of sentiment analysis on data collected from Twitter is considered to be the most commonly applied method among information retrieval work. Sentiment analysis is categorized into mainly two approaches: Machine Learning and Deep Learning. Probabilistic models and simple or complex classification methods based on learning are parts of the first approach, while Neural Networks and transformers constitute the second [23]. Examples of Machine Learning algorithms for stance detection and sentiment analysis are the Support Vector Machine (SVM) and the Naïve Bayes classifier [4]. The most used Deep Learning approaches are the Convolutional Neural Networks and Recurrent Neural Networks with its variant Long short-term memory (LSTM) and transformers such as Open-AI GPT with a unidirectional left-to-right architecture and Bidirectional Encoder Representations from Transformers (BERT) with a bidirectional approach [4].

### 1.4   Paper Outline

The rest of the study has been conducted in five sections. Section 2 represents state of the art related studies. Section 3 states some considerations regarding the Twitter API. Section 4 represents the used data set and describes the proposed methodology. In Sect. 5, the model evaluation and the experimented results have been presented. Finally, Sect. 6 includes the conclusion of the study and future directions.

## 2   Related Work

One of the first published papers regarding vaccine-related tweets applied a comprehensive Latent Dirichlet Allocation (LDA) topic modeling and dynamic sentiment representation. The research examines the subjectivity of the tweets and war-related vocabulary that creates figurative framing [22]. Also [19] calculate sentiment and subjectivity of vaccine tweets, specifically AstraZeneca, Moderna and Pfizer using TextBlob with k-nearest neighbors (k-NN) classification. Sattar et al. used a simplified text processing, TextBlob in combination with VADER to forecast the future vaccine uptake by the end of 2021, leading to mainly neutral and positive sentiments in vaccine related tweets. [18].

Mostly negative sentiments concerning vaccine tweets was the conclusion of this study [17]. In [4] four different methods were presented, namely the Bag of Words (BoW) representation, the Word embedding with machine learning, the Word embedding with Deep Learning and the BERT. By investigating the news against the results, the authors concluded that the occurrence of the tweets follows the trend of the actual events [4]. Again by using BERT model, with the different aim of vaccine misinformation detection in Twitter, this research evaluated it to be the best model of 0.98 precision test in comparison to LSTM and XGBoost [10]

A research approach that probably widens the gap between those who are skeptical against vaccine mandates and those in favor was [11] and it concludes that the anti-vaccine alert-communication aspects on Twitter should be imitated by the pro-vaccine community, by stepping into tactics of opinion manipulation. In the same journal, another interesting paper was published. They applied geographical LDA and word analysis about the US states on COVID-19 vaccines providing with analytical wordclouds and timelines of specific vocabulary used by certain groups [9]. In [12,13], the author identified twelve different communities to related topics from Twitter with uni- and bi-grams by applying graph network analysis. The dataset created from this research was published on the IEEE web portal, which is used in our paper [12,13]. In the [20] study, the authors applied the BERT model to classify the polarity and test the subjectivity of the COVID-19 tweets, with an accuracy of 93.89%.

## 3    Methodology

In this section we present a complete description of the proposed methodology. In particular, for the evaluation of the model, we used a pre-trained model in combination with a labeled dataset. COVID-19 related tweets were gathered via the Twitter API and categorized by a classifier. The classification of the collected data resulted in further analysis through time series. Figure 1 illustrates the overview of the proposed methodology.
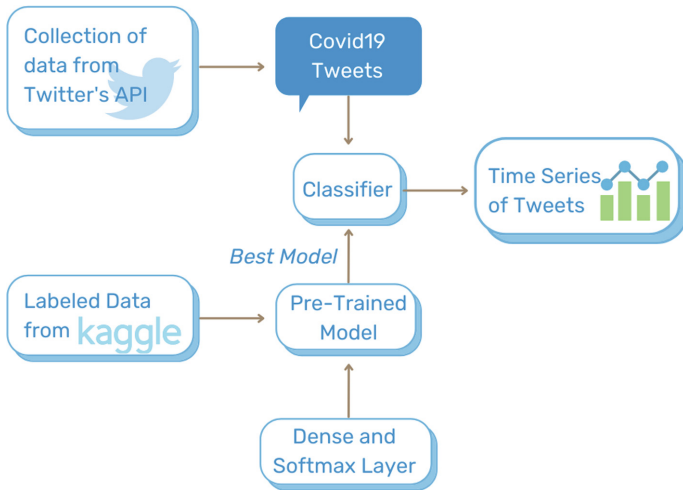


**Fig. 1.** Overview of the proposed methodology

### 3.1 Tweets Hydration

Twitter renewed its privacy policies in March 2020 in relation to the developer's platform and changed the whole tweet parsing procedure. The user profiles were not allowed to be publicly visible, shared or exposing any other sensitive information about the users. However the existence of public databases continued to be available, hence the form of IDs. These IDs can be parsed now through a procedure called hydration, which returns all the embedded information needed [13]. Of course, if the tweet is deleted by the user or filtered out by Twitter it is not hydratable.

### 3.2 Dataset

The first dataset used for this study was a publicly available dataset from a Kaggle competition. Kaggle is a platform for predictive machine learning and analytics competitions [7]. In this work, we used the Sentiment Analysis of Tweets related to the COVID-19 pandemic dataset [1] in order to validate the pre-trained model described in Sect. 3.3. The data from the source was mainly pre-processed and cleaned, but an additional data cleaning procedure has been conducted by removing duplicate entries, stopwords, digits along with hyperlinks, emojis, usernames and punctuation.
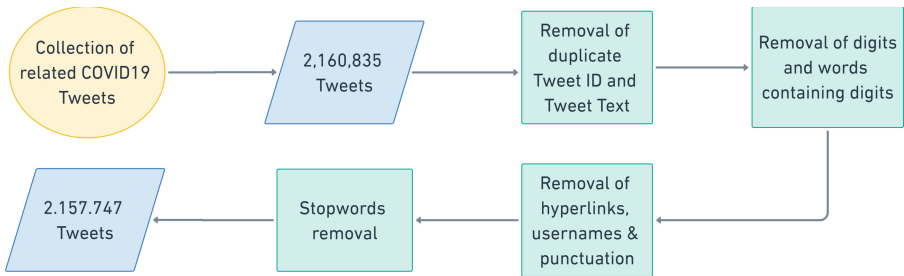


**Fig. 2.** Pre-processing procedure

The second dataset for this study contains tweets related to COVID-19, written in English, during the time period between January 01 2021 and December 31, 2021. Our search was based on a publicly accessible, large-scale COVID-19 tweets dataset [12], at a ratio approximately of 7500 tweets per day. Taking advantage of Twitter's API for hydrating, a process of retrieving a tweet's complete information using its ID. The created dataset attributes are [tweet id],[tweet creation date],[tweet text],[country],[city] and [number of retweets]. Data pre-processing is used for cleaning the collected raw data by extracting duplicate entries, stopwords, digits, as well as hyperlinks, emojis, usernames and punctuation. The pre-processing steps are shown in Fig. 2. After pre-processing, the final number of tweets was 2.157.747. Using a sample of the tweets, we mapped their location

and presented it in Fig. 3, which can be distinctive for the entire dataset. From the collected tweets, only a significantly small percentage contains location data, due to Twitter's new policy to disable the tweets' location by default. Lastly, the gathered unlabeled tweets were passed on the classifier for polarity labeling.
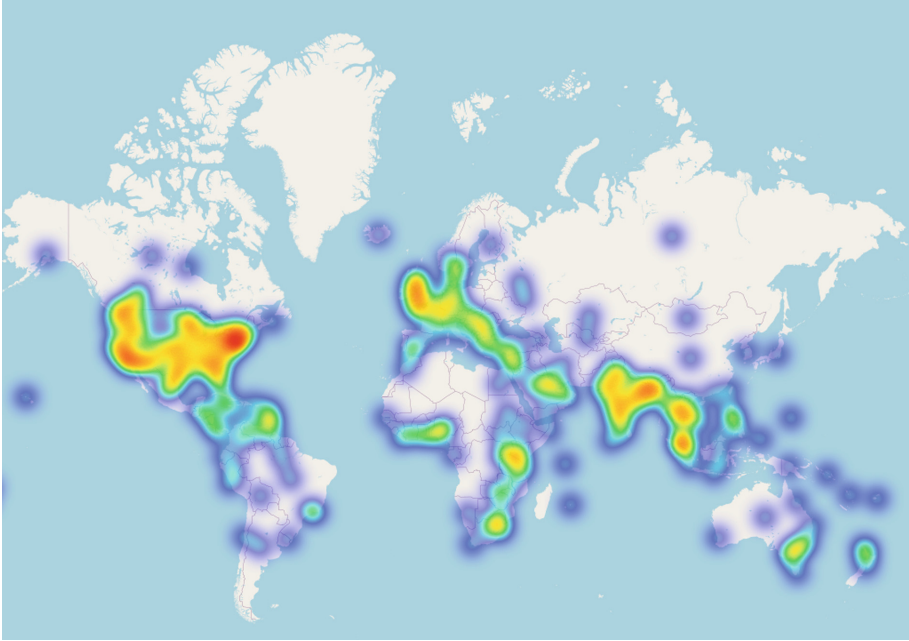


**Fig. 3.** Global tweet heat map

### 3.3   Model

To classify a tweet's polarity into the positive or the negative class, we used the BERT model [5]. The BERT model elevates standard transformers [21] by removing the unidirectionality constraint. Specifically, BERT uses a Masked Language Model (MLM) pre-training objective, which randomly masks a part of a token from the input with the main purpose to predict the original vocabulary id only by its context. The MLM objective can combine the left and the right context, providing us with the ability to pre-train a deep bi-directional transformer dissimilar to a left-to-right model. Additionally, BERT uses the *next sentence prediction* task which gives the ability to pre-train text pairs together. The model has been pre-trained on a plain unlabeled text corpus of the entire English Wikipedia and the BookCorpus. For this study, we train the whole architecture using the train set of tweets from Kaggle. Utilizing the transfer learning technique we start the training from a well known pre-trained model, instead

a random weight initialization, updating the weights of the model during the training process. On the pre-trained model we appended a dense and a softmax layer. We used the base model of BERT (Hugging Face) which uses 12 transformer encoders, 12 attention heads and 110 million parameters, instead of the large model that uses 24 transformer encoders, 16 attention heads and 340 million parameters.

## 4    Experimental Results

To investigate the model's performance, we performed 100 independent iterations. In every iteration, the Kaggle's dataset has been shuffled and split into train and test sets, with 70–30 analogy. To prevent overfitting, we performed an early trigger in our model by monitoring the train and validation loss. Once the validation loss reaches its minimum score and starts increasing, we stopped the training and saved the current model. The model showed an average accuracy of 99%, with 0.003 standard deviation, $1.510^{-5}$ variance and an average running time of 42.5 s, running on CUDA. As described, the gathered tweets have been polarized by the classifier, leading to a time series tweet average polarity per day within the year 2021.



**Fig. 4.** Tweets' polarity over the months

### 4.1    Results and Discussion

After the text classification, we measured the polarity of the average tweet per day. The polarity towards 0 presents negativity and towards 1 presents positivity. The tweets' polarity over time has been visually represented in Fig. 4. In the beginning of 2021, we observe a general neutrality lasting until the middle of March. The period between the middle of March and the end of April is characterized by a large-scale spike of positive tweets. We cross-referenced the media reports from this timespan and found that it coincided with the period with the lowest death rate the UK had seen in six months. Furthermore, at the same time most countries announced the expansion of their vaccination drive to include additional age and professional groups. The polarity decreases rapidly

in the first days of May, when, in retrospect, officials from the US authorize the vaccination of children as young as 12 years old. For a long period thereafter the tweets retain a neutral polarity, until November 2021, when a significant flow towards negativity is noticed. This can be explained by the worldwide outbreak of the Omicron variant, in November, and the subsequent take-off in the number of COVID-19 new cases.
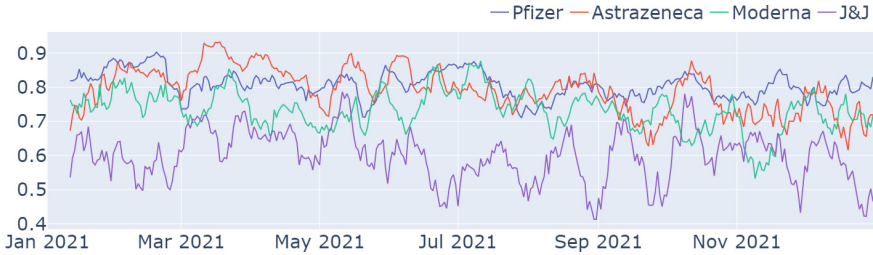


**Fig. 5.** Comparison of the COVID-19 vaccines' polarity

In addition, we isolated tweets containing the name of the four major COVID-19 vaccine companies; Pfizer/Biontech, Astrazeneca/Oxford or Vaxzevria, Moderna and Johnson & Johnson's Janssen. A visual representation of their polarity through the course of time is given in Fig. 5, where each vaccine's time series is represented in distinct colors. The results revealed that in the considered timeline, every vaccine brand has mostly positive tweets. As we can see in Table 1, the Astrazeneca/Oxford or Vaxzevria vaccine appears to be in the lead, with a score of positive tweets that reaches 82,4% of the relevant dataset.

**Table 1.** Vaccines' positive rate

| Vaccine brand | Total number of tweets | Positive tweets |
|---|---|---|
| Pfizer/Biontech | 25.447 | 81,8% |
| Astrazeneca/Oxford (Vaxzevria) | 9.474 | 82,4% |
| Moderna | 9.223 | 74,5% |
| Johnson & Johnson's Janssen | 7.297 | 61,2% |

Taking a closer look in Fig. 5, we can see that during the summer and especially within the middle of July and until the middle of August a course towards negativity has been observed. The following events took place over the summer and have been connected with the tweets' polarity concerning the vaccines:

- Israel had already offered a booster shot to its citizens.
- The FDA on July 8th stated that a booster shot was under consideration
- White House announced on August 18, the beginning of booster shots for every adult 8 months after the second dose.

The Moderna vaccine had an immediate course towards the negativity axis, in early November. Around that time the company's CEO made declarations in regard to their product's effectiveness against COVID-19.

Afterwards, we distinguished the vaccines based on their type of technology. The Pfizer/Biontech and Moderna vaccines are classified as mRNA vaccines and the Astrazeneca/Oxford and J&J vaccines as non-mRNA vaccines. Figure 6 shows the overall polarity of each vaccine type. The mRNA vaccines had overall a positive polarity over the non-mRNA vaccines. The non-mRNA vaccines show a negative peak around October 2021, which can be corresponded to the European Medicines Agency (EMA) announcement for a possible new, life-threatening side effect of J&J's vaccine, known as venous thromboembolism (VTE). Both Astrazeneca and J&J vaccines (non-mRNA) have been previously linked with a very rare thrombosis with thrombocytopenia syndrome (TTS).



**Fig. 6.** Comparison of the mRNA and Non mRNA vaccines

We concluded our study, focusing on the polarity of tweets located in the US. Figure 7 shows the average polarity of each state over the year. The states of North Dakota and Alaska have overall more positive tweets. On the other hand, the state of Oregon had tweets with the most negative polarity related to COVID-19.
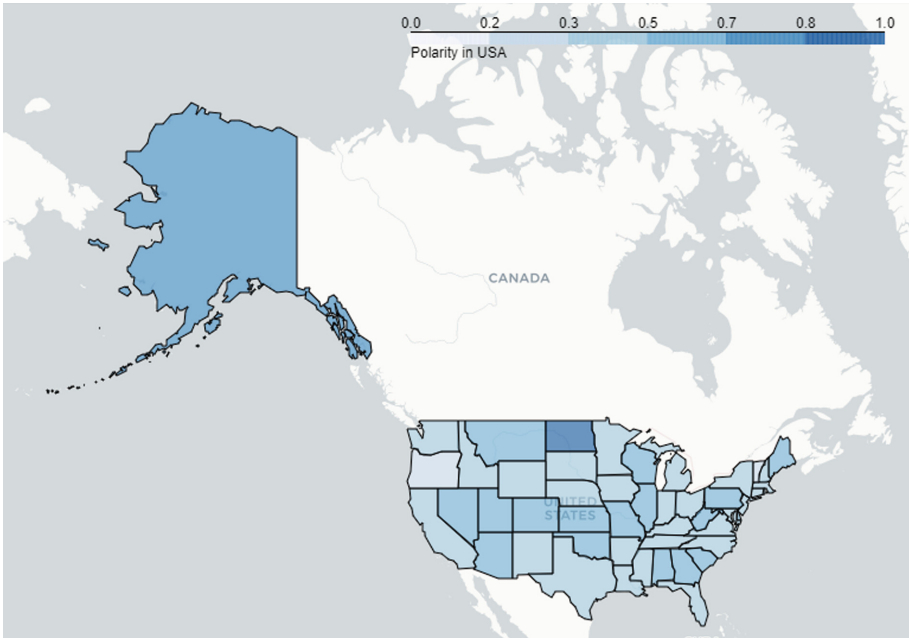
**Fig. 7.** Tweet polarity per US state

## 5   Conclusion

The primary goal of this study is to perform a sentiment analysis with reference of COVID-19 tweets within the year 2021. Taking advantage of an open labeled dataset from a Kaggle competition and a Deep Learning model, BERT, we achieved a validation accuracy as 99%. Then using the unlabeled and preprocessed tweets, we created a sentiment score for each tweet, i.e. 0 for negative and 1 for positive, which lead to an average day polarity time series. The polarity of the tweets reflects exactly the events regarding COVID-19 through the proposed period. Apart from the COVID-19 tweets, we also analyzed the polarity of the COVID-19 vaccines. In general, regardless of the vaccine's technology, the means put into effect in the fight against COVID-19 were well received from Twitter users. The effectiveness of the model brings forth new possibilities for future studies in which we can further develop and improve the performance of other text classification techniques.

# References

1. Sentiment analysis of covid-19 related tweets (2021). https://www.kaggle.com/c/sentiment-analysis-of-covid-19-related-tweets/data

2. Bullock, J., Lane, J.E., Shults, F.L.: What causes covid-19 vaccine hesitancy? ignorance and the lack of bliss in the united kingdom. Humanit. Soc. Sci. Commun. **9**(1), 1–7 (2022)

3. contributors, W.: List of covid-19 vaccine authorizations - Wikipedia, the free encyclopedia. https://en.wikipedia.org/w/index.php?title=List_of_COVID-19_vaccine_authorizations&oldid=1081403177 (2022), [Online; accessed 11-April-2022]

4. Cotfas, L.A., Delcea, C., Roxin, I., Ioanăş, C., Gherai, D.S., Tajariol, F.: The longest month: analyzing covid-19 vaccination opinions dynamics from tweets in the month following the first vaccine announcement. IEEE Access **9**, 33203–33223 (2021)

5. Devlin, J., Chang, M.W., Lee, K., Google, K., Language, A.: Bert: pre-training of deep bidirectional transformers for language understanding (2019). https://arxiv.org/pdf/1810.04805.pdf

6. Engel-Rebitzer, E., Stokes, D.C., Buttenheim, A., Purtle, J., Meisel, Z.F.: Changes in legislator vaccine-engagement on twitter before and after the arrival of the covid-19 pandemic. Hum. vaccines Immunotherapeutics **17**(9), 2868–2872 (2021)

7. Georgakopoulos, S.V., Tasoulis, S.K., Vrahatis, A.G., Plagianakos, V.P.: Convolutional neural networks for toxic comment classification. In: Proceedings of the 10th Hellenic Conference on Artificial Intelligence, pp. 1–6 (2018)

8. Gerretsen, P., et al.: Individual determinants of covid-19 vaccine hesitancy. PLOS ONE **16**(11), 1–14 (2021). https://doi.org/10.1371/journal.pone.0258462

9. Guntuku, S.C., Buttenheim, A.M., Sherman, G., Merchant, R.M.: Twitter discourse reveals geographical and temporal variation in concerns about covid-19 vaccines in the united states. Vaccine **39**(30), 4034–4038 (2021)

10. Hayawi, K., Shahriar, S., Serhani, M.A., Taleb, I., Mathew, S.S.: Anti-vax: a novel twitter dataset for covid-19 vaccine misinformation detection. Public Health **203**, 23–30 (2022)

11. Hoffman, B.L., et al.: # doctorsspeakup: lessons learned from a pro-vaccine twitter event. Vaccine **39**(19), 2684–2691 (2021)

12. Lamsal, R.: Coronavirus (covid-19) tweets dataset (2020). https://doi.org/10.21227/781w-ef42

13. Lamsal, R.: Design and analysis of a large-scale covid-19 tweets dataset. Appl. Intell. **51**(5), 2790–2804 (2021)

14. Lin, J., Ryaboy, D.: Scaling big data mining infrastructure: the twitter experience. Acm SIGKDD Explor. Newsl. **14**(2), 6–19 (2013)

15. Machado, M.D.A.V., Roberts, B., Wong, B.L.H., van Kessel, R., Mossialos, E.: The relationship between the covid-19 pandemic and vaccine hesitancy: a scoping review. Front. Public Health **9**, 1370 (2021)

16. Recio-Román, A., Recio-Menéndez, M., Román-González, M.V.: Political populism, institutional distrust and vaccination uptake: a mediation analysis. Int. J. Environ. Res. Public Health **19**(6), 3265 (2022)

17. Sarirete, A.: Sentiment analysis tracking of covid-19 vaccine through tweets. J. Ambient Intell. Humanized Comput. 1–9 (2022). https://doi.org/10.1007/s12652-022-03805-0

18. Sattar, N.S., Arifuzzaman, S.: Covid-19 vaccination awareness and aftermath: public sentiment analysis on twitter data and vaccinated population prediction in the usa. Appl. Sci. **11**(13), 6128 (2021)

19. Shamrat, F., et al.: Sentiment analysis on twitter tweets about covid-19 vaccines using NLP and supervised KNN classification algorithm. Indones. J. Electr. Eng. Comput. Sci. **23**(1), 463–470 (2021)
20. Singh, M., Jakhar, A.K., Pandey, S.: Sentiment analysis on the impact of coronavirus in social life using the BERT model. Soc. Netw. Anal. Min. **11**(1), 1–11 (2021)
21. Vaswani, A., et al.: Attention is all you need. https://arxiv.org/pdf/1706.03762v5.pdf
22. Wicke, P., Bolognesi, M.M.: Covid-19 discourse on twitter: how the topics, sentiments, subjectivity, and figurative frames changed over time. Front. Commun. **6** (2021). https://doi.org/10.3389/fcomm.2021.651997
23. Zhou, J., Ye, J.M.: Sentiment analysis in education research: a review of journal publications. Interact. Learn. Environ. 1–13 (2020)