



HAL
open science

Social Media Sentiment Analysis Related to COVID-19 Vaccines: Case Studies in English and Greek Language

Evridiki Kapoteli, Paraskevas Koukaras, Christos Tjortjis

► **To cite this version:**

Evridiki Kapoteli, Paraskevas Koukaras, Christos Tjortjis. Social Media Sentiment Analysis Related to COVID-19 Vaccines: Case Studies in English and Greek Language. 18th IFIP International Conference on Artificial Intelligence Applications and Innovations (AIAI), Jun 2022, Hersonissos, Greece. pp.360-372, 10.1007/978-3-031-08337-2_30 . hal-04668655

HAL Id: hal-04668655

<https://inria.hal.science/hal-04668655v1>

Submitted on 7 Aug 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



This document is the original author manuscript of a paper submitted to an IFIP conference proceedings or other IFIP publication by Springer Nature. As such, there may be some differences in the official published version of the paper. Such differences, if any, are usually due to reformatting during preparation for publication or minor corrections made by the author(s) during final proofreading of the publication manuscript.

Social Media Sentiment Analysis Related to COVID-19 Vaccines: Case studies in English and Greek language

Evridiki Kapoteli¹, Paraskevas Koukaras¹^[0000-0002-1183-9878], and Christos Tjortjis¹^[0000-0001-8263-9024]

The Data Mining and Research Analytics Group, School of Science and Technology,
International Hellenic University, 57001, Thessaloniki, Greece
{ekapoteli, p.koukaras, c.tjortjis}@ihu.edu.gr

Abstract. SARS-CoV-2 and its mutations are spreading around the world, threatening the human population with millions of infections and deaths. Vaccines are considered the main available weapon at hand to mitigate the spread. As a result, the development of efficient systems to understand and supervise the information dissemination, as well as the evolution of sentiments towards vaccines is critical. The goal of this research was to build and apply a supervised learning approach to monitor the dynamics of public opinion on COVID-19 vaccines using Twitter data. 1,394,535 and 61,077 tweets about COVID-19 vaccines, respectively in English and Greek, were collected, classified based on sentiment polarity and analyzed over time to gain insights into sentiment trends. Our findings reveal that overall negative, neutral, and positive sentiments were at 36.5%, 39.9%, and 23.6% in the English language dataset, respectively, whereas overall negative and non-negative sentiments were at 60.1% and 39.9% in the Greek language dataset. Policymakers and health experts could take into consideration social media sentiment analysis alongside other ways of evaluating public sentiment. Social media users are actively seeking and sharing information about pandemic-related topics, allowing governments to use social media to develop effective crisis management strategies, better inform the public with accurate and reliable news, and alleviate disease-specific concerns.

Keywords: Sentiment analysis · COVID-19 vaccines · machine learning.

1 Introduction

The first coronavirus cases emerged in mid-December 2019, but it was not until WHO declared a global pandemic in mid-March 2020 [22] that concern escalated around the world. Initially, social distance, masks, and lockdowns were the only preventive measures available to combat the pandemic, but vaccines were developed soon thereafter, and the pandemic's long-term containment depended solely on their uptake. Nevertheless, the novelty of the disease and worries regarding efficacy, safety, and vaccine development speed, as well as poor or insufficient

communication, all contributed to the population’s unwillingness to receiving the COVID-19 vaccine. The continuing spread of coronavirus and its variants necessitates the development of efficient systems to understand and monitor the flow of information and the evolution of sentiment about vaccines.

This research is motivated by the extensive use of Twitter during the COVID-19 outbreak and intends to analyze Twitter data for monitoring public opinion regarding COVID-19 vaccines. We consider the seven-month period between May 19, 2021, and November 19, 2021, to collect Twitter messages written in English and Greek. The sentiment analysis performance of various models was evaluated for each language, using an annotated dataset, and the best performing one was chosen and applied to the entire dataset. For Greek language text the sentiment analysis task was addressed as a binary task of classifying sentiment into negative and non-negative classes, and for English language text, as a three-way classification problem with negative, neutral, and positive classes.

Key contributions of our work include: i) A collection and annotation of two COVID-19 vaccination datasets, one in English and one in Modern Greek. To our knowledge, this is the first Modern Greek Twitter dataset about COVID-19 vaccines. ii) A comparative evaluation of how sentiment analysis methods work on Modern Greek, an under-resourced language for Natural Language Processing (NLP) tasks where sentiment analysis is rare, and English, the language on which most research in this area is conducted. iii) A comparison between classic machine learning models and pre-trained language models for sentiment classification. iv) An analysis of social media opinions and sentiments towards the COVID-19 vaccination among Greek individuals and the global community.

The following is an outline of the remaining paper. Section 2 provides an overview of sentiment analysis literature, with a focus on vaccines. Section 3 elaborates on the proposed research design, while Section 4 analyzes the experimental results and the trend of sentiments expressed on Twitter. The paper concludes by summarizing accomplishments, presenting limitations and future directions.

2 Background

2.1 Sentiment Analysis

Sentiment analysis focuses on emotion recognition and may employ different types of social media analytics [10], data mining and NLP methods to identify and collect information and opinions from the massive textual content available online [3]. Sentiment analysis can be applied at the document, sentence, and aspect level depending on how the text is viewed and can be classified as being machine learning-based, lexicon-based, or hybrid, depending on the techniques used. Sentiment analysis allows for real-time monitoring on all types of social media platforms [11], and thus its applications exist in nearly every field, including various types social media predictions [19] such as stock movement prediction [17, 12], election results prediction [18] and even depression detection [21]. Next

we review some recent studies with an emphasis on Twitter sentiment analysis related to vaccinations.

2.2 Sentiment Analysis for Multilingual Documents

Vaccination has always been an emotionally charged topic for societies, and as a result, a substantial amount of research has been conducted on the subject. Yuan and Crooks [23] for example, used sentiment analysis to investigate how anti- and pro-vaccination Twitter users interacted about the MMR vaccine, while Du et al. [6] suggested a hierarchical machine learning-based methodology to analyze public sentiment on HPV vaccine-related tweets. Twitter sentiment analysis related to vaccinations is a study area that has received a lot of interest, the more so because of the emergency imposed by COVID-19.

Cotfas et al. [4] compared classic Machine Learning (ML) and deep-learning algorithms to determine which one best captured public perception of the new coronavirus vaccines. The authors gathered a dataset of tweets written in English. In order to train the models, they randomly sampled a portion of the data and manually classified it as: in favor, against, or neutral to vaccination. According to their findings, classic classifiers outperformed deep-learning classifiers, with the Bidirectional Encoder Representations from Transformers (BERT) language model achieving the highest accuracy of 78.94%. When the BERT model was applied to the entire dataset, it was found that the predominant stance, either daily or entirely was neutral, whereas in favor tweets, outnumbered against tweets.

Marcec and Likic [16] performed a lexicon-based sentiment analysis and identified potential events and news that may have caused the sentiment regarding different available COVID-19 vaccines to change. During the four-month study period, the sentiment towards the Pfizer/BioNTech and Moderna vaccines has been positive and stable, while the sentiment towards the AstraZeneca/Oxford vaccine appeared to be decreasing. Another study [9] employed an open-source dataset comprising COVID-19 vaccine tweets, to determine the public stance about vaccination. The sentiment polarity of each tweet was extracted using TextBlob, and the sentiment analysis task was performed using Multinomial Naïve Bayes (MNB), Support Vector Machine (SVM), and Logistic Regression (LR) classifiers. When only positive and negative tweets were considered, the LR model performed the best, with an accuracy of 97.3%, followed by SVM and MNB, which had accuracy values of 96.26% and 88%, respectively.

Even though sentiment analysis of English-language corpora has evolved into a prominent research topic in recent years, published work on sentiment analysis of Greek-language text collections has been limited. Giatsoglou et al. [7] employed ML algorithms and deep feed-forward neural networks to classify the sentiment of online user reviews. For text representation, they trained a new model directly on a large social media corpus written in Greek, and then trained an existing language model, GreekBERT, on the same data. It was found that pre-trained language models performed better than traditional representation

models, and that they performed even better when a proper language model was trained on a smaller yet domain and task-relevant corpus.

The approach of Athanasiou and Maragoudakis [1] on sentiment analysis of Greek text, incorporated each Greek token’s translation into account as an extra input feature in the training data’s feature set, and employed the Gradient Boosting Machines algorithm. A study focusing on Twitter Sentiment Analysis applications on Greek, employed a combination of Greek lexicons and classification methods to determine how various political events influenced the emotions of Greek Twitter users in the days before the elections [2]. A corpus of annotated tweets was pre-processed to create features, and each tweet was classified using a probabilistic classifier and a hashtag-based filter.

In the context of sentiment analysis using COVID-19 data, Kydros et al. [15] collected Twitter data and used an existing lexicon enriched with specific coronavirus-related words, to assess Greek citizens’ feelings during the pandemic’s first wave in Greece. The authors found that sentiment fluctuated over time, fear dominated other emotions, and positive feelings declined while negative ones increased. The only related work in Greek on COVID-19 vaccines is the cross-sectional survey that Kourlaba et al. conducted [13], which revealed that two out of five Greek citizens were not willing or sure about getting a SARS-CoV-2 vaccine, with only 57.7% indicating they would. When the findings were compared to those of other researchers, it was found that Greeks were more hesitant to get vaccinated against COVID-19 than other Europeans.

3 Research Design

The goal of the proposed methodology is to create a model capable of predicting the sentiment of vaccine-related tweets of unknown polarity. We started by collecting two COVID-19 vaccination datasets, one with English language tweets and one with Greek language tweets, and manually annotating a subset of tweets from each dataset, that will be used to train the sentiment analysis models. The next step was to pre-process the collected tweets and determine the best representation and classification techniques. We examined Bag-of-Words and word embeddings representations with classic ML and BERT. The performance of various classifiers has been assessed using metrics such as Accuracy, Precision, Recall, and F-score, and the best-performing algorithm has been used to predict and assign a class label to each tweet. Finally, we analyzed the labeled tweets to investigate how public sentiment on Twitter has changed over time. All of the aforementioned steps are detailed in the following subsections.

3.1 Dataset Collection and Annotation

We retrieved Twitter data using several COVID-19 vaccine-related hashtags. Table 1 contains the list of hashtags that were used to collect English-language tweets. When searching for tweets in Greek, this list was modified by replacing the three first hashtags with the Greek equivalents for vaccination and vaccines.

The two distinct Twitter datasets, one for each language, were compiled between May 19, 2021, and November 19, 2021. A total of 1,394,535 English language and a total of 61,077 Greek language tweets were identified, which were reduced to 1,257,944 and 50,796 respectively after removing duplicates.

Table 1. Hashtags list for tweet search.

| | |
|----------------------------|---|
| COVID-19 Vaccination Topic | #vaccinassaveslives, #vaccinesafety, #vaccinesdontwork, #vaccine, #GetVaccinated, #CovidVaccine, #vaccinated, #vaccination, #VaccinesWork, #COVID19Vaccination, #vaxxed, #antivaxx, #vaccinationdone, #Vaccine-Deaths |
|----------------------------|---|

We then manually annotated the tweets, because we either could not identify a labeled dataset for sentiment analysis regarding COVID-19 vaccines, like in the case of Modern Greek, or identified domain-specific datasets that did not perform sufficiently well. The annotation procedure was carried out as follows.

A positive label was assigned to tweets containing expressions of support, positive attitude or emotion, and tweets describing positive situations and events. Accordingly, a negative label was assigned to tweets containing expressions of judgment, negative attitude, or emotions, as well as tweets describing negative situations and events. Tweets that did not express an emotional state were assigned a neutral label. Consequently, tweets were assigned categorical values "2", "1" or "0" indicating a positive, neutral or negative sentiment towards vaccination. However, given the low number of positive tweets present in the Greek dataset, the positive and neutral classes were merged into the non-negative one, and Greek tweets were labeled as "0" for negative and "1" for non-negative. We selected and manually annotated 2403 English and 1424 Greek language tweets, equally distributed across the classes under consideration in each case.

3.2 Data Pre-processing

The data pre-processing step follows data collection and prepares data for further analysis, whilst ensuring its quality. In this context, tweet texts were first converted to lowercase. The most frequently used emoticons were then replaced with the corresponding words, while many popular contractions, slang, and informal abbreviations were also replaced with their original forms. Elements like URLs, username mentions, numbers, and special characters were discarded since they did not contribute to our analysis. Regarding hashtags, we deleted the hashtag symbol ("#") while keeping the content because they are frequently used instead of normal words and contain valuable information. For Greek language tweets, we considered an additional step at the start of the cleaning process that includes replacing accented vowels with unaccented ones.

In addition, we eliminated stopwords depending on the classification model being tested, and conducted lemmatization of the remaining words. Lemmatiza-

tion is the process of eliminating a word’s inflectional endings and returning it to its base form. Stopwords are words that appear frequently in a corpus, but do not provide additional meaning in an analysis, and need to be eliminated. In our case, negative terms like "no" and "don’t", which are commonly found in stopwords lists express sentiment, and eliminating them would completely change the stance of the text. As a result, we used a stopwords list provided by the Natural Language Toolkit (NLTK) library [20] and a list of Greek stopwords [8] after eliminating negative keywords.

3.3 Text Representation and Classification

Using the annotated dataset as training data, the performance of the following prominent classifiers was reviewed in our research: Support Vector Machine (SVM), Logistic Regression (LR), Random Forest (RF), and Extreme Gradient Boosting (XGBoost). However, before using any classification algorithm, text data needs to be translated into a numerical feature vector. The simplest method is known as Term Frequency - Inverse Document Frequency (TF-IDF) and weights each term and assigns it a TF and IDF score, where the former determines the number a term appears in a document, and the latter indicates how common or rare a term is across the entire collection of documents. The second most common method is word embeddings. This study employs the Word2Vec algorithm, which builds a vocabulary by considering the words that occur in a collection of sentences more times than a numeric threshold specified by the user. It then implements the continuous bag-of-words model or the skip-gram model to learn a D-dimensional vector representation of each word from the input documents. In addition to learning individual word representations, using Word2Vec there is the option of learning word context representations with BERT. BERT is a transformer-based language model, that can be fine-tuned to a user’s specifications using its pre-training as a base layer. In this research, we selected BERT_{BASE}¹ as described by the research team [5] and GreekBERT², the Greek equivalent of BERT_{BASE} [14].

4 Results

First, we used a TF-IDF Vectorizer with its n-gram parameter set to (1,2), referring to unigrams and bigrams. Then, we trained a Word2Vec model on our collected training data to obtain vector representations for all the unique words present in the corpus, using the skip-gram training algorithm and setting the dimensionality of its word vectors equal to 100. Lastly, we further trained BERT_{BASE} and GreekBERT language models, on our collected English and Greek language social media corpus respectively, fine-tuning them using a batch size of 16, a learning rate of 2e-5 and four epochs.

¹ <https://huggingface.co/bert-base-uncased>

² <https://huggingface.co/nlpauieb/bert-base-greek-uncased-v1>

As already mentioned, we evaluated the selected models by considering if the elimination of stopwords and the implementation of lemmatization improved their performance. We present the findings of the scenario that most classifiers perform best in. For the TF-IDF scheme, this entails performing both stopwords elimination and lemmatization in each tweet. For the English language Word2Vec model, we only performed stopwords removal, while for the Greek language Word2Vec model we did not apply stopwords removal or lemmatization. Fig. 1 lists the hyperparameters optimized for each tunable algorithm using the grid search technique, as well as their optimal values in the best models.

| | | SVM | LR | RF | XGBoost |
|---------|----------|--|--------------------------|--|---|
| English | TF-IDF | 'C': 2, 'gamma': 'scale', 'kernel': 'linear' | 'C': 10, 'max_iter': 100 | 'max_features': 'log2', 'min_samples_split': 4, 'n_estimators': 200 | 'gamma': 0.1, 'max_depth': 3, 'n_estimators': 200 |
| | Word2Vec | 'C': 10, 'gamma': 0.1, 'kernel': 'rbf' | 'C': 5, 'max_iter': 500 | 'max_features': 'auto', 'min_samples_split': 2, 'n_estimators': 300 | 'gamma': 0.01, 'max_depth': 9, 'n_estimators': 200 |
| Greek | TF-IDF | 'C': 1, 'gamma': 'scale', 'kernel': 'linear' | 'C': 10, 'max_iter': 100 | 'max_features': 'log2', 'min_samples_split': 2, 'n_estimators': 500 | 'gamma': 0.1, 'max_depth': 3, 'n_estimators': 100 |
| | Word2Vec | 'C': 5, 'gamma': 0.5, 'kernel': 'rbf' | 'C': 10, 'max_iter': 500 | 'max_features': 'auto', 'min_samples_split': 5, 'n_estimators': 500 | 'gamma': 0.1, 'max_depth': 9, 'n_estimators': 100 |

Fig. 1. Hyperparameters and their optimal values in each case.

Next, Fig. 2 shows the accuracy score achieved by BERT and each ML algorithm. When employing the TF-IDF technique for representing text in English, we can see that Logistic Regression was the best performing classifier, followed by SVM, with 85% and 84% accuracy scores, respectively. The results indicate that the use of word embeddings improved the performance of all algorithms except LR for which the accuracy score did not change. The RF and XGBoost classifiers showed the greatest improvement in accuracy, whereas SVM outperformed all models, and was the best classifier when word embeddings representations were used. Having an accuracy of 91% the BERT_{BASE} model exceeded the performance of all the aforementioned classifiers. We also observed that for both weighting schemes, TF-IDF and word embeddings, the considered algorithms performed worse in terms of precision, recall, and f-score, for the neutral class, which did not hold for the BERT language model.

When examining the classification performance results for the Greek language dataset, we observe that SVM was the best performing classifier for the TF-IDF weighting method, with an accuracy of 86%, followed by LR. The other two algorithms performed relatively poorly with accuracy rates of 78% and 76%. When it comes to word embeddings, the performance of RF and XGBoost classifiers improved by 4% and 3% respectively, whilst the accuracy rates of the SVM and LR models were slightly reduced. Despite the decrease in its accuracy, the SVM model still outperformed the other classifiers. However, the best performance among the explored approaches was achieved when fine-tuning the GreekBERT language model with an accuracy of 93%. This model also outperformed the others in terms of precision, recall, and f-score across all classes, negative and non-negative.

| Representation Technique & Classifier | TF-IDF | | | | Word2Vec | | | | BERT |
|--|------------|-----------|-----------|----------------|------------|-----------|-----------|----------------|------|
| | <i>SVM</i> | <i>LR</i> | <i>RF</i> | <i>XGBoost</i> | <i>SVM</i> | <i>LR</i> | <i>RF</i> | <i>XGBoost</i> | |
| English | 0.84 | 0.85 | 0.79 | 0.77 | 0.86 | 0.85 | 0.82 | 0.83 | 0.91 |
| Greek | 0.86 | 0.85 | 0.78 | 0.76 | 0.83 | 0.81 | 0.82 | 0.79 | 0.93 |

Fig. 2. Classification Accuracy scores for English and Greek language models.

To summarize, most of the trained classifiers performed well, however, the two fine-tuned BERT models had the highest accuracy scores. Consequently, after discarding the training tweets we applied the fine-tuned BERT_{BASE} model to all 1,255,554 tweets in the English language corpus, and the fine-tuned GreekBERT model to all 49,375 tweets in the Greek language corpus. The objective of this process was to determine the polarity of each tweet, so we could analyze the classified tweets over time and identify sentiment changes towards the vaccination topic, which could occur in response to vaccine-related news or events.

4.1 Trend Analysis

Fig. 3 shows the percentage of tweets distributed across the different sentiment categories in both languages. Based on the distribution of English language tweets, the "Neutral" category had the highest percentage of tweets, reaching almost 40%, while the "Positive" category with 23.6% was outnumbered by the "Negative" category with 36.5%. On the other hand, the overwhelming majority of Greek language tweets (60.1%) fell into the "Negative" category, with the remaining 39.9% falling into the "Non-negative" category.

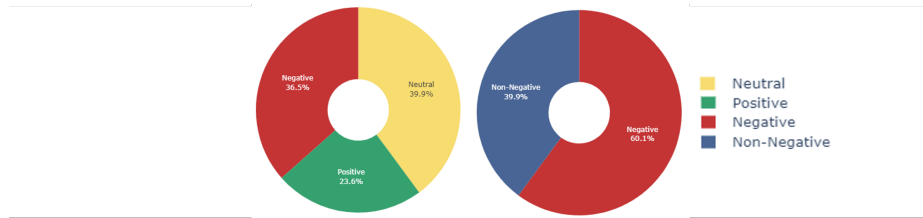


Fig. 3. Tweets percentage distribution by sentiment category collectively for English (left) and Greek (right).

Since these are aggregated results and may conceal fluctuations that occurred within the time range, we further plotted the daily timeline showing the evolution of vaccine-related tweets based on sentiment, for the study period. The daily progression of negative and non-negative sentiments on the left side of Fig. 4 shows that the dominant sentiment was negative, but its trend is closely followed by the non-negative trend. The time series for the English language dataset, presented on the right side of Fig. 4, shows that the prevailing sentiment remained neutral until July 15, 2021, when negative sentiment started to dominate, with a few exceptions. The number of positive tweets did not exhibit dramatic fluctuations except for September 17, 2021, and October 21, 2021, when it peaked.

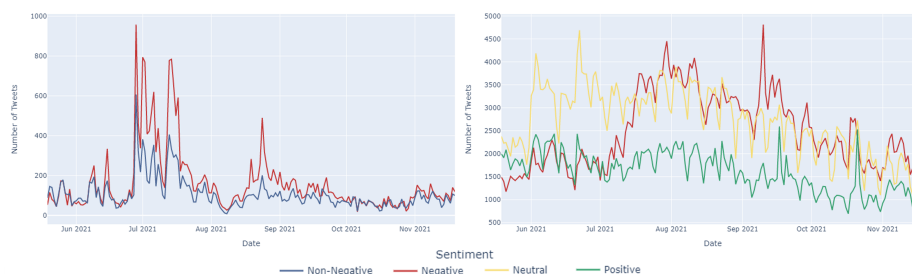


Fig. 4. Time series between May 19, 2021 and November 19, 2021, for the collected tweets on a daily basis based on sentiment for Greek (left) and English (right).

Next, we studied word clouds displaying the one hundred most frequently occurring words for each sentiment class, that were generated using the log-likelihood values. According to Fig. 5, the words "death", "unvaccinated", "covidiot" and "antivaxxer" appeared frequently in negative English language tweets, indicating that many Twitter users expressed their displeasure with people who refused to receive the vaccine. People were still concerned about the mask mandates and the coronavirus repercussions, as evidenced by terms like "mask", "mandate", "sick" and "risk." Words that appeared often in neutral tweets such as "vaccine", "dose", "available", and "date", indicate a need for information on the vaccination program. Interestingly, there were several keywords related to India, such as "indiafightscorona", "pune" and "covaxin", implying that a significant number of tweets originated from Indian users. Finally, terms like "vaccinated", "thank", and "great" appeared in the positive word cloud, indicating that people support and receive COVID-19 vaccines, as well as encourage others to do so by using hashtags like "getvaccinated" and "vaccineswork". Some more words shown in a smaller size are "happy", "grateful", "protect" and the hashtag "staysafe".

Moreover, according to Fig. 5 and in the case of Greek language word clouds, the most interesting conclusions can be drawn from the negative one. It is evident that it was dominated by references to Greece's prime minister and the government in general, the mass media, as well as rude and inelegant characterizations of them. Other negative words refer to vaccine deniers, vaccine side-effects, coronavirus variants and deaths, whereas some of the most popular keywords associated with non-negative sentiment referred to the pandemic, covid cases, and mutations, vaccines, dose, children, and hashtags like "covidgreece", and "rollingupsleeves".

5 Conclusion

This paper follows the evolution of sentiment towards COVID-19 vaccines between May 19, 2021, and November 19, 2021, by creating a ML based sentiment analysis model using pre-annotated tweets, that is capable of classifying Twitter posts written in English or Greek. Several established classifiers and language models were examined for both languages, with their accuracy rates ranging from



Fig. 5. Word clouds for the English (top) and Greek (bottom) datasets, showing the most frequent appearing words in vaccine-related tweets for each sentiment category.

76% to 93%. The proposed framework classifies English language tweets into positive, neutral, and negative categories employing BERT with 91% accuracy. Greek language tweets are classified as negative and non-negative, employing GreekBERT with an accuracy of 93%.

When the English dataset was studied collectively, the prevailing sentiment was neutral, but daily, neutral was only dominant during the first months, as the sentiment shifted to negative in the months that followed, with the exception of a few days. When we compare the percentage of tweets belonging into each sentiment category on the beginning and end dates of the study period, we find that negative sentiment increased by nearly 12%, positive sentiment decreased by 11%, and the percentage of neutral sentiment remained almost stable. Except for a few days negative sentiment dominated the Greek dataset both overall and daily. When the percentage of tweets on May 19 and November 19 were compared, it was found that negative tweets grew by 17% while non-negative tweets declined by 17%.

As the vaccination process is still hampered by several barriers, and new cases and deaths are growing worryingly, our work demonstrates that social media sentiment analysis can yield useful insights, which governments and health experts can use to develop effective crisis management strategies, better inform the public and plan ahead of time to prevent disease-specific concerns. Quick responses and actions facilitated by social media analysis, aimed at minimizing and preventing negative emotional and psychological impacts will enhance global health and well-being amid crises such as the SARS-CoV-2 pandemic.

5.1 Limitations and Future Work

Although Twitter is a valuable data source for studying real-time social media content about coronavirus vaccines, its users are not representative of the general Greek- and English-speaking public, and their tweets simply reflect netizens' views and emotions about vaccination. Another limitation is linked to data collection, since we queried Twitter using limited sets of vaccine-related hashtags,

which may have been incomplete. Next, determining the class to which each tweet belongs introduces subjectivity, since each opinion may have different interpretations. We should also keep in mind that social media data contains a lot of noise, which is difficult to completely eliminate, and that the selection of data processing techniques is subjective and may significantly impact the research's outcomes. Additional limitations stem from the use of ML, which can process large amounts of data considerably faster than human approaches, but still struggles to detect irony and sarcasm in tweets.

For that reason, the primary focus of future research should be on building better-performing sentiment classification models and conducting the analysis with a larger dataset, acquired from multiple social network sites or sources other than social media. It would be also interesting to detect and classify emotions in tweets, such as happiness, fear, trust, which has the potential to improve citizens' and society's understanding. Another future direction involves studying the geographic distribution of tweets and their sentiment, comparing how sentiment changes across different countries around the world.

References

1. Athanasiou, V., Maragoudakis, M.: A novel, gradient boosting framework for sentiment analysis in languages where nlp resources are not plentiful: A case study for modern greek. *Algorithms* **10** (3 2017). <https://doi.org/10.3390/A10010034>
2. Belevelis, D., Tjortjis, C., Psaradelis, D., Nikoglou, D.: A hybrid method for sentiment analysis of election related tweets. In: 2019 4th South-East Europe Design Automation, Computer Engineering, Computer Networks and Social Media Conference (SEEDA-CECNSM). pp. 1–6 (2019). <https://doi.org/10.1109/SEEDA-CECNSM.2019.8908289>
3. Cambria, E., Schuller, B., Xia, Y., Havasi, C.: New avenues in opinion mining and sentiment analysis. *IEEE Intelligent Systems* **28**, 15–21 (2013). <https://doi.org/10.1109/MIS.2013.30>
4. Cofas, L.A., Delcea, C., Roxin, I., Ioanăș, 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). <https://doi.org/10.1109/ACCESS.2021.3059821>, <https://ieeexplore.ieee.org/document/9354776>
5. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR abs/1810.04805* (2018), <http://arxiv.org/abs/1810.04805>
6. Du, J., Xu, J., Song, H., Liu, X., Tao, C.: Optimization on machine learning based approaches for sentiment analysis on hpv vaccines related tweets. *Journal of Biomedical Semantics* **8**, 9 (12 2017). <https://doi.org/10.1186/s13326-017-0120-6>
7. Giatsoglou, M., Vozalis, M.G., Diamantaras, K., Vakali, A., Sarigiannidis, G., Chatzisavvas, K.C.: Sentiment analysis leveraging emotions and word embeddings. *Expert Systems with Applications* **69**, 214–224 (3 2017). <https://doi.org/10.1016/J.ESWA.2016.10.043>
8. Greek stopwords collection. <https://github.com/stopwords-iso/stopwords-el>, accessed: 2022-02-23

9. Khakharia, A., Shah, V., Gupta, P.: Sentiment analysis of covid-19 vaccine tweets using machine learning. *SSRN Electronic Journal* (2021). <https://doi.org/10.2139/ssrn.3869531>
10. Koukaras, P., Tjortjis, C.: Social media analytics, types and methodology. In: *Machine Learning Paradigms*, pp. 401–427. Springer (2019)
11. Koukaras, P., Tjortjis, C., Rousidis, D.: Social media types: introducing a data driven taxonomy. *Computing* **102**(1), 295–340 (2020)
12. Koukaras, P., Tsihli, V., Tjortjis, C.: Predicting stock market movements with social media and machine learning. In: *Proceedings of the 17th International Conference on Web Information Systems and Technologies - Volume 1: WEBIST*, pp. 436–443. INSTICC, SciTePress (2021). <https://doi.org/10.5220/0010712600003058>
13. Kourlaba, G., Kourkouni, E., Maistrelis, S., Tsopela, C.G., Molocha, N.M., Triantafyllou, C., Koniordou, M., Kopsidas, I., Chorianopoulou, E., Maroudimanta, S., Filippou, D., Zaoutis, T.E.: Willingness of greek general population to get a covid-19 vaccine. *Global Health Research and Policy* **6** (12 2021). <https://doi.org/10.1186/s41256-021-00188-1>
14. Koutsikakis, J., Chalkidis, I., Malakasiotis, P., Androutsopoulos, I.: GREEKBERT: the greeks visiting sesame street. *CoRR* **abs/2008.12014** (2020), <https://arxiv.org/abs/2008.12014>
15. Kydros, D., Argyropoulou, M., Vrana, V.: A content and sentiment analysis of greek tweets during the pandemic. *Sustainability* **13**(11) (2021). <https://doi.org/10.3390/su13116150>, <https://www.mdpi.com/2071-1050/13/11/6150>
16. Marcec, R., Likic, R.: Using twitter for sentiment analysis towards astrazeneca/oxford, pfizer/biontech and moderna covid-19 vaccines. *Postgraduate Medical Journal* (8 2021). <https://doi.org/10.1136/postgradmedj-2021-140685>
17. Nousi, C., Tjortjis, C.: A methodology for stock movement prediction using sentiment analysis on twitter and stocktwits data. In: *2021 6th South-East Europe Design Automation, Computer Engineering, Computer Networks and Social Media Conference (SEEDA-CECNSM)*. pp. 1–7 (2021). <https://doi.org/10.1109/SEEDA-CECNSM53056.2021.9566242>
18. Oikonomou, L., Tjortjis, C.: A method for predicting the winner of the usa presidential elections using data extracted from twitter. In: *2018 South-Eastern European Design Automation, Computer Engineering, Computer Networks and Society Media Conference (SEEDA-CECNSM)*. pp. 1–8 (2018). <https://doi.org/10.23919/SEEDA-CECNSM.2018.8544919>
19. Rousidis, D., Koukaras, P., Tjortjis, C.: Social media prediction: a literature review. *Multimedia Tools and Applications* **79**(9), 6279–6311 (2020)
20. Steven, B., Ewan, K., LoperEdward: *Natural Language Processing With Python*. O'Reilly (2009)
21. Wang, X., Zhang, C., Ji, Y., Sun, L., Wu, L., Bao, Z.: A depression detection model based on sentiment analysis in micro-blog social network. In: Li, J., Cao, L., Wang, C., Tan, K.C., Liu, B., Pei, J., Tseng, V.S. (eds.) *Trends and Applications in Knowledge Discovery and Data Mining*. pp. 201–213. Springer Berlin Heidelberg, Berlin, Heidelberg (2013)
22. Timeline: Who's covid-19 response. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/interactive-timeline#event-72>, accessed: 2022-02-23
23. Yuan, X., Crooks, A.T.: Examining online vaccination discussion and communities in twitter. *ACM International Conference Proceeding Series* pp. 197–206 (7 2018). <https://doi.org/10.1145/3217804.3217912>