27th International Conference on Knowledge-Based and Intelligent Information & Engineering Systems (KES 2023)

# SART & COVIDSentiRo: Datasets for Sentiment Analysis Applied to Analyzing COVID-19 Vaccination Perception in Romanian Tweets

Alexandra Ciobotaru, Liviu P. Dinu

*Faculty of Mathematics and Computer Science, University of Bucharest, Academiei Street 14, Bucharest, 010014, Romania*

## Abstract

Vaccination is an important subject of discussion adjacent to the COVID-19 pandemic. Sentiments generated online by this topic are worth analyzing using opinion mining tools, and it is interesting to do so in online content written in an under-researched language, like Romanian. For this reason, we modified and enlarged an existing sentiment analysis dataset comprised of Romanian tweets labeled as negative or positive. The resulting dataset, SART (Sentiment Analysis from Romanian Tweets), comprised of three classes (positive, negative, and neutral) containing 1300 Romanian tweets each, was used to train two different sentiment analysis models: a fastText-based one and a fine-tuned BERT model. We further show the usefulness of the sentiment analysis model by analyzing the sentiment of Romanian tweets regarding vaccination using a corpus created and collected by the authors between January 2021 and February 2022 (COVIDSentiRo).

*Keywords:* Type your keywords here, separated by semicolons ;

## 1. Introduction

The increasing quantity of data in online makes opinion mining an extremely useful tool. People like to share their opinions [1] and find comfort in talking to each other over social media about the topics that are of most significance to them. The Twitter platform can be a useful tool for opinion mining, as the platform is stable over time, thus making it easy to extract data with scrappers and analyze people's opinions. Although generally, Twitter is not one of Romanians' favorite social media platforms, since the COVID-19 pandemic started, the number of Romanian Twitter

---

* Corresponding author. Tel.: +40-740-016-235
  *E-mail address:* alexandra.ciobotaru@unibuc.ro

users has increased by approximately 100,000, as stated by statista.com. As shown by Florea and Roman [2], most Romanian Twitter users are between 30 and 35 years old and residents of Bucharest city. The linguistic style of tweets is usually informal, with abbreviations, idioms, and jargon [3], and the maximum length of a tweet recently increased to 280 characters, making these texts suitable for vectorization with most LLM (Large Language Model) tokenizers.

For analyzing Romanian Twitter users' opinions regarding vaccination against COVID-19 disease, a sentiment analysis model is needed. There are already some sentiment analysis datasets in Romanian. One is LaRoSeDa [4], which contains Romanian reviews gathered from popular e-commerce sites, 7500 positive and 7500 negative. Another sentiment analysis dataset contains Romanian movie reviews, roughly 11000 positive and 7000 negative, and is available online. But the only dataset for sentiment analysis from tweets in Romanian, to our knowledge, is the one described by Istrati and Ciobotaru [5], which is a dataset created for brand analysis, containing roughly 1000 positive tweets and 1000 negative tweets.

As the task of opinion mining is greatly dependent on the way the sentiment analysis model is created, in view of the fact that one cannot use a sentiment analysis model created on movie reviews or product reviews to analyze tweets, we opted for Istrati and Ciobotaru dataset [5] for analysing vaccination opinions in Romanian, but we modified it by reconsidering its correctness and enlarged it by adding texts to each class and also by adding a neutral class, as it has been shown that neutral examples help obtain superior classification results [6]. We name this dataset SART (Sentiment Analysis from Romanian Tweets) and release it to the public.

The sentiment analysis model created based on the SART dataset is further used to infer vaccination-related texts extracted in the time frame 01.01.2021-01.02.2022, an operation that led to the creation of the first corpus of COVID-19 related tweets in Romanian. We name this corpus COVIDSentiRo and release it to the public as well, considering that the pandemic was an important event that affected our society in ways that deserve further insights by the research community.

The inference also revealed interesting correlations between the weekly number of people vaccinated against COVID-19 and the sentiment of tweets, be it positive, negative, or neutral. These results are presented in Section 5.

The main contributions that this work makes are the following:

1. We enhance and enlarge an existing sentiment analysis dataset for Twitter by correcting annotation mistakes, adding more tweets, and also adding a neutral class, and we introduce it as the **SART dataset**, useful for analysing the sentiments in Romanian tweets.
2. Based on SART, we create a sentiment analysis model, which we use to analyse Romanian tweets regarding vaccination against COVID-19 on the novel **COVIDSentiRo** corpus, scrapped from Twitter in the time frame 01.01.2021-01.02.2022.
3. We present the correlative results between Romanians' sentiments on Twitter about vaccination against COVID-19 and the weekly number of vaccinated Romanians.
4. The datasets, sentiment analysis model, and results are freely available at `https://github.com/Alegzandra/KES-2023`.

## 2. Related Works

In terms of related works about the COVID-19 pandemic opinion mining, there are mainly two approaches: one that uses unsupervised clustering and another based on supervised techniques that make use of sentiment analysis machine learning models.

For the Romanian language, regarding the first approach, Reveiu and Arghir [9] use a clustering technique to analyse Twitter users opinions regarding the pandemic crisis in terms of brand analysis for the main retailers in Romania. Likewise, Harba et al. [10] analyse perceptions in the fine-dining restaurant industry in Bucharest, Romania, using customer reviews before and during the pandemic.

Regarding the English language, Syed at al. [11] use a H-TF-IDF approach to extract terms from English tweets and analyse public opinion regarding COVID-19 in the month of January 2020. They use a hybrid approach by clustering tweets and using classical machine learning models like SVM and Logistic Regression to get trends.

Another similar analysis using the English language is described by Zhang et al. in their work [12]. The authors analyse the opinions of Twitter users over a period of time between February and October 2020 in some cities in Canada and the United States and find that, in general, the users' sentiments regarding COVID-19 are positive regarding masks but negative regarding the topic of vaccination.

Cioban and Vîntoiu [13] compute the frequency of words in Romanian Reddit posts and find correlations between fluctuations of extremely high and extremely low sentiment scores and the declared COVID-19 cases.

In order to analyse vaccine hesitancy, Lanyi et al. [14] gather 91,473 tweets between November 2020 and August 2021, and do cluster analysis according to topic and sentiment.

Another interesting analysis is the one described by Yousefinaghani et al. [15] in their work, where the authors gather 4,552,652 publicly available tweets posted between January 2020 and January 2021, and find their polarity using a lexicon-based method of analysis described in [16]. They also indicate trends in the general public's attitudes toward vaccination throughout the course of the study period.

Naseem et al. [17] create a large-scale sentiment dataset in English, COVIDSENTI, which contains 90,000 COVID-19-related tweets collected from February to March 2020, labeled as positive, negative, and neutral. Their analysis shows that Twitter users had positive opinions during the stay at home order in February, but their opinions shifted the next month. Further, Jalil at al. [18] apply pre-processing and classification techniques to the COVIDSENTI dataset, and their experiments show an accuracy of 96.66% on the task of sentiment analysis on COVID-19 related tweets.

Other than English and Romanian, a study worth mentioning is the one shown by La Gatta et al. [19], which investigates how Italians perceived the pandemic of COVID-19 by applying emotion analysis over Italian tweets related to this subject from January 2020 to February 2022.

To our knowledge, there is no study that analyses Romanian opinions regarding the COVID-19 pandemic using machine learning methods. Inspired by the works presented, we aspire to fill this gap by creating such an analysis. During this process, we enhance and release to the public a general dataset for sentiment analysis in Romanian tweets, SART (Sentiment Analysis in Romanian Language), which we used to create a sentiment analysis model using machine learning techniques. We further apply the model by predicting the sentiment of 19,319 tweets regarding the COVID-19 pandemic and vaccination (COVIDSentiRo), a corpus extracted using specific keywords, in the time frame January 2021-February 2022. Our analysis shows cues on how the vaccination campaign was received by the population, and we hope that SART and COVIDSentiRO will continue to be used in future analyses in terms of painting the right picture on the impact of pandemic events on the Romanian population.

## 3. SART

### 3.1. Upgrading a Sentiment Analysis Dataset

Istrati and Ciobotaru [5], using their custom dataset, created several classical machine learning models for sentiment analysis and also two more modern models, one that implied fine-tuning the Romanian BERT [20] and also a fastText approach [21]. What we noticed was that the fine-tuned BERT model created was trained using classical preprocessing of the custom dataset, comprised of, among others: stemming, removing punctuation, diacritics, lowercasing, translating emoticons, and removing stopwords, with an additional BERT tokenization layer on top as well. The same classical vectorization process was applied for their fastText model as well, although in the fastText library description it is stated clearly that such preprocessing steps are not necessary.

For this reason, we retrained the fine-tuned BERT model on the custom dataset for sentiment analysis, using the same train/test split proportions, with the only preprocessing consisting mainly in anonymisation: eliminating usernames using regex and eliminating proper names and organization names using Romanian name entity recognition [22], as we did not want our model to be biased towards certain users, persons, or organizations. Lastly, we eliminated hyperlinks, newlines, tabs, redundant spaces, and other artefacts, using regex.

For fine-tuning the Romanian BERT we used the same set-up described in [5]. We created a simple classifier that loads the weights from the romanian-bert-cased and adds on top a linear layer, activated using the Softmax function. We trained the model for 10 epochs using the cross-entropy loss function and the AdamW optimizer for the learning rate. As expected, at the end of the training process we obtained an accuracy of 88% on the test set, which is **eight** percent higher than the accuracy obtained by the authors in [5] using double tokenization.

As the texts in the dataset created and described in the mentioned article were related only to specific brands, which were used as key words when extracting the data from Twitter, we considered it opportune to broaden their dataset by adding more general texts, together with their corresponding label, positive or negative, and a neutral class as well.

First, we manually verified the dataset, which was preprocessed for anonymity as described above. From the tweets contained in that dataset, we deleted those that did not evoke the labeled sentiment, in our opinion, and also those tweets that did not make any more sense after masking usernames, proper nouns, and organizations.

Further, we extracted tweets in Romanian from Twitter in two stages: first bulk of tweets was extracted on the 5th of February 2022, and the second bulk on the 20th of February 2022. The sentiment of the extracted tweets was predicted using the fine-tuned BERT model for sentiment analysis created above. We kept and verified only the texts classified with a probability of over 90%, and during the verification step, we deleted the tweets that were not conveying the labelled sentiment.

Our general annotation approach corresponds to the annotation guidelines described in [23] which involve two annotating stages: a first round of labeling made by the first annotator and a second round of verification made by the second annotator. The texts over which the two annotators did not agree were deleted. It is important to note that the original custom dataset was labeled by two annotators, but in a sequential fashion, meaning that each text was verified by only one annotator.

We also analyzed the tweets classified with probabilities lower than 50%, correctly assuming that these texts are most probably neither positive nor negative. After manually checking the selected texts, we chose approximately 500 tweets as neutral.

In this way, we collected 500 positive, 500 negative, and 500 neutral tweets.

The final step in creating our enhanced dataset for sentiment analysis was combining the tweets that remained from the custom original dataset, curated after verification, with the ones collected in the steps described above. The final resulting dataset consists of 1300 negative tweets, 1300 positive tweets, and 1300 neutral tweets. The neutral tweets were obtained by combining the 500 neutral tweets from inference with 800 neutral tweets extracted randomly from the neutral class of the REDv1 dataset [24]. We chose this dataset because it is the only single-label dataset containing neutral Romanian tweets available open-source.

### 3.2. Creating the Sentiment Analysis Model

Based on the data obtained, the next step was to create the sentiment analysis model. We created two models, one based on the Romanian BERT model, cased version [20], and one based on fastText word embeddings [21].

We shuffled the data and split it into 3120 labelled texts for training, 390 for validation and 390 for testing, according to the 80%-10%-10% scheme.

#### 3.2.1. BERT Based Model

We fine-tune the Romanian BERT, cased version, using the same set-up described in Subsection 3.1. Results are shown in Table 1. Metrics were computed using the "classification_report" method from the scikit-learn Python library [25]. It can be seen that the classification of negative and positive tweets has the highest precision, while classification of the neutral tweets have the highest recall. Thus, the model returns more relevant results than irrelevant ones for the negative and positive classes, and it returns most of the relevant results, regardless of their irrelevancy, for the neutral class. Also, the averaged precision, recall, and f-score are almost equal to each other, with the precision being 85% and the recall and f-score equaling 84%.

Table 1. Classification report for both the BERT-based model and the fastText-based model

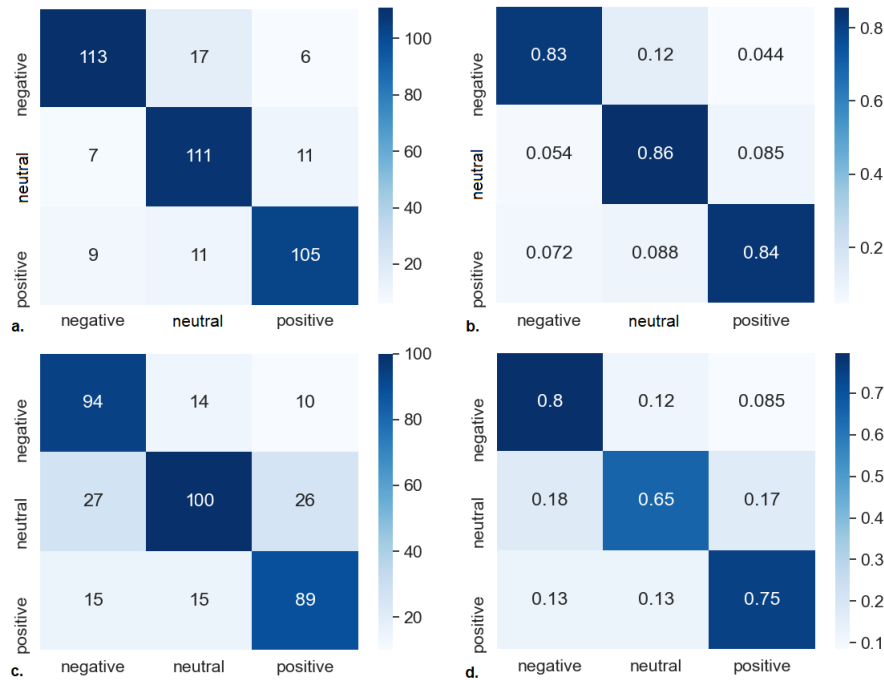| | BERT-based | | | fastText-based | | | |
|---|---|---|---|---|---|---|---|
| **Sentiment** | **Prec.** | **Rec.** | **F-score** | **Prec.** | **Rec.** | **F-score** | **Supp.** |
| Negative | **0.88** | 0.83 | 0.85 | **0.80** | 0.69 | 0.74 | 136 |
| Neutral | 0.80 | **0.86** | 0.83 | 0.65 | **0.78** | 0.71 | 129 |
| Positive | **0.86** | 0.84 | 0.85 | **0.75** | 0.71 | 0.73 | 125 |
| Macro avg | 0.85 | 0.84 | **0.84** | 0.73 | 0.73 | 0.73 | 390 |

Fig. 1. Confusion matrices for the created sentiment analysis models: a. BERT-based - Classic; b. BERT-based - Normalized; c. fastText-based - Clasic; d. fastText-based - Normalized

In Figure 1 a. and b., we show confusion matrices for the BERT-based model, both normalized and classic. Normalization is useful for observing the class with the highest accuracy in the model. The highest miss-classified value in Figure 1 a. is 17, meaning that 17 negative tweets were confused as being neutral, but the support is also higher for this label (136 tweets), meaning that it is natural to find more miss-classifications for this category, than the others. Also, in Figure 1 a., we can compute the overall accuracy of the model by adding the values on the diagonal and dividing by the whole number of texts analyzed in the test set. The accuracy of the fine-tuned BERT model is 84%.

From Figure 1 b., it can be seen that the model overall performs best when classifying the neutral class (0.86), but the difference from the least well-classified label, the negative one, is very small - only 0.03. So it can be said that overall, the model performs in a balanced way across all three classes.

### 3.2.2. FastText Based Model

For the fastText-based model [21], on the same training data and using the same data for validation and testing as for the BERT-based model described in Subsection 3.2.1, we used the method "train_supervised" from the fasttext Python library, which applies the supervised text classification method described by Joulin et al. in [26]. In a nutshell, the architecture of the model is similar to the CBOW described by Mikolov et al. in [27], but an important particularity of the FastText model is that it vectorizes the texts by creating character-level n-grams, which are further fed to a linear classifier. For optimization, the model uses SGD (Stochastic Gradient Descent) and a learning rate that decays linearly, and for computing the probabilities on the output classes, it uses Hierarchical Softmax.

Using the "autotuneDuration" parameter of the "train_supervised" method, we have let our model autotune itself on the validation file for a duration of 10 minutes. The results are shown in Table 1.

It can be seen that both negative and positive classes have the highest precision, while the highest recall is reached by the neutral class. Overall, all metrics are smaller than what the BERT-based model achieved. The averaged precision, recall, and f-score are equal to the accuracy of the model, which is 73%.

Studying Figure 1 c. and d., where we show confusion matrices for the fastText based model, both normalized and classic, we can see that the model classifies best the negative tweets, and worst the neutrals, while the highest number

of miss-classifications in the classical confusion matrix (1 c.) is of 27 neutral tweets, which were miss-classified by the fastText-based model as negative.

### 3.2.3. Error Analysis for the BERT-based model

As the BERT-based model described in Subsection 3.2.1 had better performances, we decided to use this model for inference in Section 5, but not without a previous error analysis. Analysing the texts that were miss-classified by the model, we noticed that the majority of them contained sarcasm in varying degrees, although we should be aware that sarcasm can occasionally be challenging for human readers as well.

## 4. COVIDSentiRo

In order to gather the data needed for our use-case, we scrapped tweets in Romanian in the time-frame 01.01.2021-28.02.2022, using *snscrape* Python library. This time frame was chosen to cover the beginning of the vaccination campaign in Romania as well as the period after the booster vaccination campaign ended.

The query words used for scrapping, along with their English translations, are shown in Table 2.

Table 2. Query words used for scrapping, with English translations.

| Lang | Query Words |
| --- | --- |
| ro | vaccin, vaccinuri, vaccinare, vaccinează, vaccinări, vaccinat, vaccinată, vaccinați, vaccinate |
| en | vaccine, vaccines, vaccination, vaccinate, vaccinations, vaccinated, vaccinated, vaccinated, vaccinated |
| ro | vaccinații, vaccinatele, nevaccinat, nevaccinați, nevaccinate |
| en | vaccinated, vaccinated, unvaccinated, unvaccinated, unvaccinated |
| ro & en | vax, antivax, anti-vax, anti-vaxx, antivaxx |
| ro & en | moderna, pfizer, biontech, astrazeneca |
| ro & en | covid, covid-19, coronavirus, corona |

After the scrapping process was finished, a total of 128,853 tweets resulted, but after some operations applied to the resulting dataframe, only 34437 tweets remained. The operations were to remove duplicate tweets and eliminate all tweets containing hyperlinks, as most of those were news regarding the analyzed topic.

Further, we used roner [22] to anonimize the texts by masking entities of type PERSON with the sequence "< |PERSON| >", like described in [29]. We used regex patterns to mask usernames with "< |USERNAME| >", emails with "< |EMAIL| >", telephone numbers with "< |TEL| >". In order to prepare the texts for inference, we further eliminated redundant spaces and artefacts from scrapping, as well as the sequence masks.

Next, using *langdetect* Python library, we eliminated tweets that were not in Romanian, as we found a significant number of them, although the snscrape library was used with 'lang=ro' parameter. As langdetect is based on a probabilistic model, we noticed that it classifies erroneously short texts, so we eliminated texts shorter than 30 characters. By doing this, we eliminated roughly 1,000 tweets.

After all these steps were done, 19,319 tweets remained in the final dataset, ready for sentiment analysis.

## 5. Correlation Results

The tweets prepared for inference in Section 4 were further classified using the fine-tuned BERT model created in Section 3. This step resulted in a corpus containing 10,434 negative, 7,944 neutral, and 941 positive tweets.

In Figure 2 we represent the number of vaccinated people per week, alongside the number of tweets labeled as positive, neutral, or negative, in the period of time analyzed. For creating the graph representing the number of vaccinated people per week, we downloaded the weekly number of COVID-19 vaccinations, using the COVID19 R package [30], from which we aggregated the weekly number of vaccinated people in Romania.

There can be observed two major peaks in vaccination: one between the 1st week of 2021 and week 25 of 2021 (*first vaccination phase*) and another between week 39 of 2021 and week 52 of 2021 (*second vaccination phase*). In
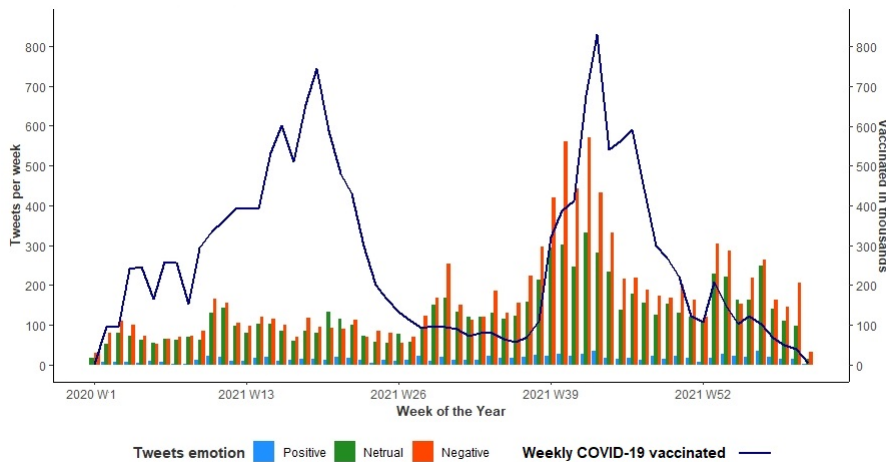
Fig. 2. Number of vaccinated people per week, and the number of weekly tweets labelled by sentiment.

Table 3. The correlation coefficient (with its corresponding p-value) between weekly tweets by sentiment and weekly COVID-19 vaccinations (p=0.05).

| Sent | Wks 1-21 | Wks 41-62 | Wks 1-62 |
|------|----------|-----------|----------|
| Neg  | 0.41 (0.06) | **0.64 (<p)** | **0.28 (<p)** |
| Neu  | **0.51 (<p)** | **0.62 (<p)** | **0.3 (<p)** |
| Pos  | **0.53 (<p)** | 0.34 (0.13) | 0.19 (0.13) |

the first 21 weeks, there was a growing trend in vaccination among people, and then the trend became descendant until the beginning of the second vaccination campaign in Romania, which included the administration of the booster dose and the apparition of the very contagious Delta wave.

Also, starting with the $20^{th}$ of September 2021, the COVID-19 certificate became compulsory in Romania. This date belongs to week **38** in the analyzed period. Looking at Figure 2 it can be observed that this week is the starting point for the second vaccination phase in Romania, which corresponds as well to an increase in debate on the Twitter platform regarding this subject. While in the first vaccination phase this subject was not discussed so much on Twitter, starting with the second vaccination phase, the prevailing sentiments of Twitter users regarding this subject are mostly negative or neutral.

In Table 3 we computed correlation scores between the number of tweets per week by sentiment and the number of COVID-19 **vaccinations**, for three periods of time: *weeks 1-21* (first vaccination phase), *weeks 41-62* (second vaccination phase), and *the whole time-frame*. Correlations were computed using *cor.test* R function with the Pearson method [31]. A correlation value is considered significant statistically only if its corresponding p-value is lower than 0.05. Thus, in both tables, the emphasized values are statistically significant.

Looking at Table 3, we can see that over the entire time period (Weeks 1-62), there are no correlations large enough to be considered. Broken down by time periods, in the first vaccination phase there is a correlation between the weekly number of positive and neutral tweets and the number of vaccinated people, while in the second vaccination phase there are correlations between negative and neutral tweets and the number of vaccinated people.

In general, the sentiments found in the COVIDSentiRO dataset express attitudes towards vaccination and not towards the disease itself. Thus, the positive tweets found in the first vaccination phase are mainly jokes regarding vaccination, about choosing a certain type of vaccine over another (Moderna, Pfizer, or AstraZeneca) and there can be seen a general happiness towards the fact that a vaccine for COVID-19 finally exists. In the second vaccination phase, there can be seen negative tweets addressed towards the antivaccinists, but there can also be found negative tweets written by antivaccinists towards those who are pro-vaccine. These arguments frequently center on issues like the requirement for the green certificate and the negative side effects of vaccines. However, in order to generalize these results, a detailed analysis is needed in future studies.

## 6. Conclusions

We enhance a sentiment analysis dataset with Romanian tweets by rechecking annotations, adding new annotated tweets, and adding a neutral class as well. We use this new dataset (SART) to train a sentiment analysis model. The best model performances are achieved by fine-tuning the Romanian BERT model. Further, we collect Romanian tweets regarding COVID-19 over a time frame that covers the two major vaccination phases against COVID-19 that took place in Romania. The corpus used for analysis is called COVIDSentiRo and is the first corpus created for analyzing opinions regarding COVID-19 in Romanian tweets.

We show the usefulness of our sentiment analysis model by predicting and analysing the sentiment of each text in COVIDSentiRO. Results show a correlation between positive tweets and the number of vaccinated people in the first 21 weeks of the period analyzed, corresponding to the first vaccination phase, and a correlation between negative tweets and the last 21 weeks of the period analyzed, corresponding to the second vaccination phase, a time when vaccination against COVID-19 became compulsory in Romania as the "green certificate". This suggests that the general opinion on Twitter regarding the COVID-19 vaccination subject shifted from positive in the first vaccination phase to negative towards the end of the second vaccination phase, a period of time that also generated more debate among the Twitter users, as noticeable by a general increase in the number of tweets.

## References

[1] Wlodarczak, P., Ally, M. & Soar, J. Opinion Mining in Social Big Data. *SSRN*. (2015,2)
[2] Florea, A. & Roman, M., The Profile Of Social Media Users In Romania: Individual Characteristics And The Number Of Social Connections. (2019,5)
[3] Martínez-Cámara, Eugenio and Martín-Valdivia, Maria and López, L. and Montejo-Ráez, Arturo, Sentiment analysis in Twitter. *Natural Language Engineering*, vol 20 (2014)
[4] Tache, A., Mihaela, G. & Ionescu, R. Clustering Word Embeddings with Self-Organizing Maps. Application on LaRoSeDa - A Large Romanian Sentiment Data Set. *Proceedings Of The 16th Conference Of The European Chapter Of The Association For Computational Linguistics: Main Volume*. pp. 949-956 (2021,4), https://aclanthology.org/2021.eacl-main.81
[5] Istrati, L. & Ciobotaru, A. Automatic Monitoring and Analysis of Brands Using Data Extracted from Twitter in Romanian. *Intelligent Systems And Applications*. pp. 55-75 (2022)
[6] Koppel, M. & Schler, J. Using Neutral Examples for Learning Polarity. *IJCAI-05, Proceedings Of The Nineteenth International Joint Conference On Artificial Intelligence, Edinburgh, Scotland, UK, July 30 - August 5, 2005*. pp. 1616-1617 (2005), http://ijcai.org/Proceedings/05/Papers/post-0312.pdf
[7] Alexandra Ciobotaru, November 26, 2022, "SART", IEEE Dataport, doi: https://dx.doi.org/10.21227/5fnc-tk84
[8] Alexandra Ciobotaru, November 26, 2022, "COVIDSentiRO", IEEE Dataport, doi: https://dx.doi.org/10.21227/a7az-aw35
[9] Reveiu, A. & Arghir, D., Mining Social Media To Identify The Immediate Impact Of Covid-19 Pandemic On The Romanian Retailers: Early Findings. *New Trends In Sustainable Business And Consumption*. **1225** (2020)
[10] Harba, J., Tigu, G. & Davidescu, A. Exploring Consumer Emotions in Pre-Pandemic and Pandemic Times. A Sentiment Analysis of Perceptions in the Fine-Dining Restaurant Industry in Bucharest, Romania. *International Journal Of Environmental Research And Public Health*. **18**, 13300 (2021)
[11] Syed, M., Decoupes, R., Arsevska, E., Roche, M. & Teisseire, M. Spatial Opinion Mining from COVID-19 Twitter Data. *International Journal Of Infectious Diseases*. **116** pp. S27 (2022), https://www.sciencedirect.com/science/article/pii/S1201971221009577, Abstracts from the Eighth International Meeting on Emerging Diseases and Surveillance, IMED 2021
[12] Zhang, Q., Yi, G., Chen, L. & He, W. Text mining and sentiment analysis of COVID-19 tweets. *CoRR*. (2021), https://arxiv.org/abs/2106.15354
[13] Cioban, Ș. & Vîntoiu, D. The rebellious social network reaction to COVID-19. *Studia Universitatis Babes-Bolyai Sociologia*. **65**, 111-130 (2020)
[14] Lanyi, K., Green, R., Craig, D. & Marshall, C. COVID-19 Vaccine Hesitancy: Analysing Twitter to Identify Barriers to Vaccination in a Low Uptake Region of the UK. *Frontiers In Digital Health*. **3** (2021)
[15] Yousefinaghani, S., Dara, R., Mubareka, S., Papadopoulos, A. & Sharif, S. An analysis of COVID-19 vaccine sentiments and opinions on Twitter. *International Journal Of Infectious Diseases*. **108** pp. 256-262 (2021), https://www.sciencedirect.com/science/article/pii/S1201971221004628
[16] Hutto, C. & Gilbert, E. VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. *Proceedings Of The International AAAI Conference On Web And Social Media*. **8**, 216-225 (2014,5), https://ojs.aaai.org/index.php/ICWSM/article/view/14550
[17] Naseem, U., Razzak, I., Khushi, M., Eklund, P. & Kim, J. COVIDSenti: A Large-Scale Benchmark Twitter Data Set for COVID-19 Sentiment Analysis. *IEEE Transactions On Computational Social Systems*. **8**, 1003-1015 (2021)
[18] Jalil, Z., Abbasi, A., Javed, A., Badruddin Khan, M., Abul Hasanat, M., Malik, K. & Saudagar, A. COVID-19 Related Sentiment Analysis Using State-of-the-Art Machine Learning and Deep Learning Techniques. *Frontiers In Public Health*. **9** (2022), https://www.frontiersin.org/articles/10.3389/fpubh.2021.812735

[19] Valerio La Gatta and Vincenzo Moscato and Marco Postiglione and Giancarlo Sperli, A Longitudinal Study on Italian Reactions to the Different Narratives of Covid-19. *IEEE Intelligent Systems*. (2023)

[20] Dumitrescu, S., Avram, A. & Pyysalo, S. The birth of Romanian BERT. *Findings Of The Association For Computational Linguistics: EMNLP 2020*. pp. 4324-4328 (2020,11), https://aclanthology.org/2020.findings-emnlp.387

[21] Bojanowski, P., Grave, E., Joulin, A. & Mikolov, T. Enriching Word Vectors with Subword Information. *arXiv preprint arXiv:1607.04606*. (2016)

[22] Dumitrescu, S. & Avram, A. Introducing RONEC - the Romanian Named Entity Corpus. *Proceedings Of The 12th Language Resources And Evaluation Conference*. pp. 4436-4443 (2020,5), https://aclanthology.org/2020.lrec-1.546

[23] Orbach, M., Toledo-Ronen, O., Spector, A., Aharonov, R., Katz, Y. & Slonim, N. YASO: A Targeted Sentiment Analysis Evaluation Dataset for Open-Domain Reviews. *Proceedings Of The 2021 Conference On Empirical Methods In Natural Language Processing*. pp. 9154-9173 (2021,11), https://aclanthology.org/2021.emnlp-main.721

[24] Ciobotaru, A. & Dinu, L. RED: A Novel Dataset for Romanian Emotion Detection from Tweets. *Proceedings Of The International Conference On Recent Advances In Natural Language Processing (RANLP 2021)*. pp. 296-305 (2021,9), https://ranlp.org/ranlp2021/proceedings.pdf

[25] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. & Duchesnay, E. Scikit-learn: Machine Learning in Python. *Journal Of Machine Learning Research*. **12** pp. 2825-2830 (2011)

[26] Joulin, A., and Grave, E., Bojanowski, P. & Mikolov, T. Bag of Tricks for Efficient Text Classification *arXiv preprint arXiv:1607.01759*. (2016)

[27] Mikolov, T., Chen, K., Corrado, G. & Dean, J. Efficient estimation of word representations in vector space *arXiv preprint arXiv:1301.3781*. (2013)

[28] Barbieri, F. & Saggion, H. Automatic Detection of Irony and Humour in Twitter. *ICCC*. (2014)

[29] Ciobotaru, A., Constantinescu, M., Dinu, L. & Dumitrescu, S. RED v2: Enhancing RED Dataset for Multi-Label Emotion Detection. *Proceedings Of The 13th Language Resources And Evaluation Conference (LREC 2022)*. pp. 1392-1399 (2022), http://www.lrec-conf.org/proceedings/lrec2022/pdf/2022.lrec-1.149.pdf

[30] Guidotti, E. & Ardia, D. COVID-19 Data Hub. *Journal Of Open Source Software*. **5**, 2376 (2020)

[31] Best, D. & Roberts, D. Algorithm AS 89: The Upper Tail Probabilities of Spearman's Rho. *Journal Of The Royal Statistical Society. Series C (Applied Statistics)*. **24**, pp. 377-379 (1975)

## 7. Comment

In this second version of the article we addressed the requirements of the reviewers by making the following changes:

- We have added a paragraph at the end of Section 5, explaining how to understand the positive and the negative content in the context of the disease.
- We fixed spelling errors.
- We have added the author's names near references in Related Works Section.
- We rewrote the fine-tuned BERT model description by adding more details upon the structure of the modelin Subchapter 3.1.
- We have added an extended description of the fasttext model in Subsection 3.2.2, along with 2 new references: [26] and [27].
- We have added the names of the authors of this article and their affiliations.