

# Bilingual emotion analysis in social media throughout the COVID19 pandemic in Portugal

**Alina Trifan**  
DETI/IEETA  
University of Aveiro  
Portugal  
alina.trifan@ua.pt

**Sérgio Matos**  
DETI/IEETA  
University of Aveiro  
Portugal  
aleixomatos@ua.pt

**Pedro Morgado**  
School of Medicine  
University of Minho  
Portugal  
pedromorgado@med.uminho.pt

**José Luís Oliveira**  
DETI/IEETA  
University of Aveiro  
Portugal  
jlo@ua.pt

## Abstract

This paper presents preliminary work on the topic of emotion analysis on Twitter, in the context of the coronavirus pandemic in Portugal. We collected, curated and analyzed covid-related tweets of users in Portugal posted during four and a half months (January to May 2020) in order to understand the evolution of the six basic emotions reflected in these tweets. We analyzed tweets written in both English and Portuguese. In this first step of our work we correlate this information with key events of the evolution of the pandemic in Portugal during March, which was the most critical period in Portugal. Our findings show that the sentiment analysis of covid-related tweets is consistent with our hypothesis that negative emotions would intensify as the pandemic progressed. The preliminary results obtained stand as proof of concept that the analysis of real-time tweets or other social media messages through sentiment analysis can be an important tool for behavioural and well-being tracking.

## 1 Introduction

The focus of attention of health care providers around the world for the last several months has been the problem of the new coronavirus (COVID-19) and its spread. In addition to efforts at various levels to prevent the spread of the disease and other worrisome conditions, special attention should be paid to mental health and care. According to similar epidemics and pandemics, in such cases, serious concerns such as fear of death can arise among patients, and feelings of loneliness and insecurity can develop. Moreover, people who are quarantined lose face-to-face connections and traditional social interventions, which significantly lowers personal and mental well-being.

We present in this short paper preliminary work on the analysis of social media posts for inferring

the prevalence and evolution of the six basic human emotions in Portugal, during the first four and a half months of the pandemic. Our hypothesis is that negative emotions, and consequently a decrease in personal well-being and possibly mental health, would be predominant and that their intensity would follow the evolution of the COVID-19 cases in Portugal. We are interested in understanding this evolution and performing sentiment analysis as an initial step of a broader goal of monitoring mental health and well-being among social networks users in these challenging times. This paper is structured in five more sections. We shortly discuss next the current background in sentiment analysis, mental health and well-being monitoring in the context of COVID-19 pandemic. In Section 3 we detail the process of data collection and curation. We present our analysis methodology in Section 4 and the results obtained so far in Section 5. Finally, we conclude the paper and discuss future research steps in Section 6.

## 2 Background

In the current context, global attention has largely been focused on the infected patients and the front-line responders, with some marginalised populations in society having been overlooked (Ho et al., 2020). Previous research has revealed a profound and broad spectrum of psychological impact that outbreaks can inflict on people (Hall et al., 2008; Müller, 2014).

Several recent publications in the area of mental health try to raise alarm flags with respect to the importance of finding new methods for responding to mental health needs during the pandemic time. Xiang et al. (Xiang et al., 2020) claim that the mental health needs of patients with confirmed COVID-19, patients with suspected infection, quarantined family members, and medical personnel

have been poorly handled throughout this time. Similarly, Yao et al. (Yao et al., 2020) express their concerns with regards to the effect of the epidemic on people with mental health disorders. The obvious solution for Wind et al. (Wind et al., 2020) to continue mental health care within a pandemic is to provide mental health care at a warm distance by video-conferencing psychotherapy and internet interventions. This idea is further supported by a study among Chinese citizens (Gao et al., 2020), in which Gao et al. assess the prevalence of mental health problems and examine their association with social media exposure. Their findings highlight that there is a link between the two, which suggests that the government need to pay more attention to mental health problems and their online manifestation. A portuguese study also found that individuals previously receiving psychotherapeutic support benefit if they did not interrupt the process as a consequence of the outbreak (Moreira et al., 2020).

The already available scientific evidence previously introduced strongly suggests that a shift in mental health care provision towards online prevention and treatment in the near future is needed. As such, the widespread use of social media combined with the rapid development of computational infrastructures to support big data, and the maturation of natural language processing and machine learning technologies offer exciting possibilities for the improvement of both population-level and individual-level mental health (Conway and O'Connor, 2016).

### 3 Data collection

This work builds on prior research contributions, both national and international, that enabled the collection of both social media and clinical statistics data. For social media data, we collected four and a half months of Twitter<sup>1</sup> posts whose content contained coronavirus related vocabulary, as explained next. The data collection process followed the pipeline recently published by Chen et. al (Chen et al., 2020). According to Twitter's Terms and Conditions, Chen et. al have released Tweet IDs, which are unique identifiers tied to specific tweets containing coronavirus related terms. The collection of ids is available on GitHub<sup>2</sup>.

Based on this dataset, we queried the Twitter API and obtained the complete dataset. For each

tweet in the collection we retrieved tweet content (text, URLs, hashtags) and authors' metadata. Following authors' suggestions, we used Hydrator<sup>3</sup>, Twitter's search API<sup>4</sup> and Twarc<sup>5</sup> for retrieving the data. All tweets included in this collection contain coronavirus related vocabulary, both as hashtags or mentions/keywords in the text of the tweet. We collected tweets posted from 21st of January to 31st of May, 2020. More details on the collection process can be found in the original publication (Chen et al., 2020).

We filtered the collected tweets based on the tweet's geo-location or user's location. When users post tweets from a GPS-enabled device or they tag a location in their tweet, this information can be retrieved as longitude and latitudine coordinates. Unfortunately, only roughly 1–3% of Twitter messages are geocoded (Paul and Dredze, 2017). Because the number of geo-located tweets in the corpora we collected is indeed limited, we used the location defined on the user's profile for extracting tweets of Portuguese users. Using Tweepy<sup>6</sup> we filter for user locations that match a dictionary of city names in Portugal. We filter out false positives by removing locations that also match names of cities or states in other countries, namely Brazil (e.g. a match for Porto Alegre, a city in Brazil, replaces a match for the Portuguese city Porto). We ended up with a list of 76898 distinct users and 117772 distinct tweets. This list of ids will be made publicly available upon paper acceptance.

For clinical statistics and the identification of key events of the manifestation of the COVID-19 pandemic in Portugal we relied on an open-source repository maintained by Data Science for Social Good Portugal<sup>7</sup>, an open community of data scientists that tackle relevant and current societal issues. This repository contains daily updates of the statistic information released by the Portuguese Ministry of Health regarding COVID-19 cases in Portugal.

### 4 Data analysis methods

Psychologist Paul Eckman identified six basic emotions that he suggested were universally experienced in all human cultures. These emotions were

<sup>3</sup><https://github.com/DocNow/hydrator>

<sup>4</sup><https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets>

<sup>5</sup><https://github.com/DocNow/twarc>

<sup>6</sup><https://www.tweepy.org/>

<sup>7</sup><https://github.com/dssg-pt/covid19pt-data>

<sup>1</sup>[www.twitter.com](http://www.twitter.com)

<sup>2</sup><https://github.com/echen102/COVID-19-TweetIDs>

happiness (joy), sadness, disgust, fear, surprise, and anger (Ekman, 1992). Happiness is often classified as a positive emotion, while the remaining five basic emotions are classified as negative. We studied the basic emotions reflected in Portuguese user tweets at the word level, using Natural Language Processing (NLP), both in English and in Portuguese. Words with basic discrete emotions were tabulated, counted and their frequency was measured.

We first filtered tweets by language and addressed only tweets written in English or Portuguese. In the preprocessing step we performed tokenization, stemming and lemmatization, lowercase conversion and stopwords removal for each of the two languages. To this purpose we chose Python<sup>8</sup> as a programming language and the Natural Language Toolkit<sup>9</sup> as NLP framework. For the tweets written in English, we used Empath (Fast et al., 2016), a lexicon mined from modern text on the web, using a combination of deep learning and crowdsourcing, to retrieve the counts and frequencies of the six basic emotions. Each of these emotions corresponds to an Empath lexical category, which is formed by the emotion word and a collection of other similar words that convey the same emotion.

For the tweets written in Portuguese, a psychiatrist member of our research team defined a semantic cluster of Portuguese words that are associated to each of the six basic emotions. These semantic clusters are presented in Table 1. For both languages, we used the lexicon words’ frequency in a tweet to measure the strength of a specific emotion. For example, if two words in a tweet were coded as joy or two occurrences of terms belonging to the joy semantic cluster, and one word was coded as fear, joy was counted twice and fear was counted once. These counts were then normalized with respect to the total number of tokens in a tweet.

## 5 Preliminary results

We hereby present preliminary results on the sentiment analysis of the evolution of the six basic emotion of Portuguese Twitter users, along with other statistical information regarding the corpus that we curated. We consider this curation of tweets written by Portuguese users an important scientific contribution by itself. We will publish the list of

<sup>8</sup>[www.python.org](http://www.python.org)

<sup>9</sup><https://www.nltk.org/>

Emotion	Semantic cluster
Joy (alegria)	feliz; felicidade; espetacular; esperança; expectativa; fantástico; wow; alegria
Sadness (tristeza)	triste; deprimido; depresso; tristeza
Disgust (nojo)	nojo; contaminação; repulsa; contágio
Fear (medo)	medo; ansioso; preocupado; apreensivo; nervoso;
Surprise (surpresa)	surpreendido; surpresa; inesperado
Anger (raiva)	wtf; merda; pqp; fdx; revoltado; zangado; irritado; enervado

Table 1: Portuguese language semantic clusters defined for each of the six basic emotions.

Month	# PT	# EN	# Other
January	2679	2262	79
February	17277	3312	1175
March	29158	11663	1745
April	10022	6401	1006
May	12407	5733	941

Table 2: Number of tweets per month, global numbers and tweet counts discriminated by language - Portuguese (PT), English (EN) and other languages. It is important to note that in January we only retrieved posts from the 21st of January to 31st of January.

tweets ids, annotated with the language in which the tweet was written, as open source upon paper acceptance. We want, on one hand, to encourage further exploration of this data for gaining general well-being insights in Portugal and on the other hand, to serve as possible comparison base for similar studies going on in other countries.

Table 2 overviews the number of tweets posted by Portuguese users, by month, from 21st of January to 31st of May. For each month, we indicate the number of tweets written in Portuguese and in English, as well as the total number of tweets written in other languages. Among the tweets written in other languages than Portuguese and English, Spanish, French and Italian were the most predominant ones.

Just by looking at the total number of tweets per month we can see that March was the most critical month of the pandemic in Portugal. The

first COVID-19 positive cases were confirmed by the end of February but in a small number. Partial lockdown started in Portugal by mid-March and it slowly progressed into total lockdown by the end of the month. The number of tweets related to COVID-19 published in this period speaks for the increased interest and concern that this pandemic had during the month of March. Another important remark is that the ratio between the number of tweets written in Portuguese and the number of tweets written in English during February and March is superior to the ones in April and May, which might suggest an increased focus on the development of the national COVID-19 situation, rather than worldwide.

The main hypothesis of the present study is that the tweets collected would reflect more negative emotions throughout the pandemic, up to its peak, which in Portugal it is estimated to have happened around 23rd-25th of March 2020. The evolution of COVID-19 cases in Portugal, both in terms of infected population and number of suspected cases is presented in Fig. 1.

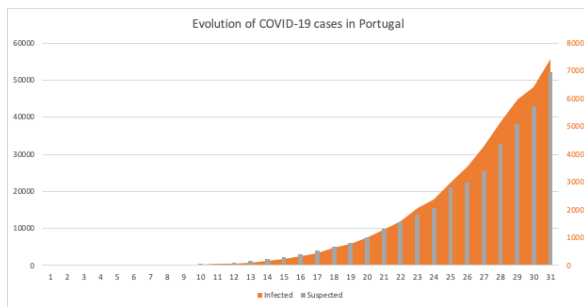


Figure 1: COVID-19 cases in Portugal throughout the month of March.

We examined the patterns of the tweets and the evolution of the basic emotions to see whether they were consistent with our expectations. We performed this analysis for all the tweets collected so far. However, due to the late-breaking nature of this research, we will focus here on the month of March, the most critical one so far in Portugal.

We present the number of tweets per day written by Portuguese users during the month of March in Fig. 2. Figure 3 shows the normalized evolution of the six basic emotions in covid-related tweets written in English by Portuguese users over the month of March 2020. As we can see, the negative emotions are predominant with little joy or positive sentiment present in these tweets. Moreover, it is important to note the spike in fear around the

period when home isolation started. On the 7th of March the first COVID-19 cases were found among student population and in the following day 2 universities were closed.

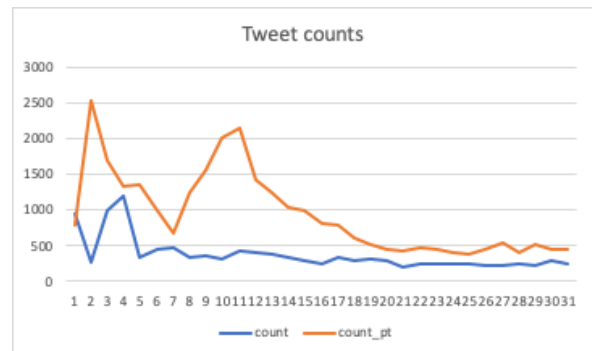


Figure 2: Number of tweets per day during March. Tweets written in English in blue and tweets written in Portuguese in orange. There are two clear spikes in the number of Portuguese tweets in the beginning of the month, with the first COVID-19 suspected cases and a second spike around the period when home isolation was decreed in Portugal.

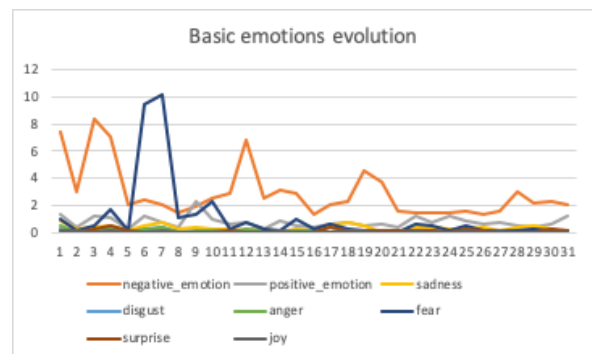


Figure 3: Evolution of the basic emotion frequencies in tweets written in English and Portuguese during March. We also include the evolution as categories of *negative\_emotion* and *positive\_emotion*.

## 6 Conclusions and future directions

Social media reflects the lives of a population and their attitudes and as such, social media sentimental analysis and behavioural tracking has the potential of becoming an important tool into turning public health provision more personalized and efficient. We presented in this paper the evolution of the six basic emotions over Twitter, during the COVID-19 pandemic in Portugal, which can be correlated with the disease evolution. We found the emotional patterns of the tweets largely consistent with our expectations, with more negative tweets posted in the beginning of the pandemic, up to its peak.

As future work, we are interested in analyzing this dataset into more detail in order to leverage mental health and well-being status knowledge of the social media users. As such, we will explore psycholinguistic features and previously trained models on social media corpora of mental-health related issues in order to better understand the impact of this pandemic.

## Acknowledgments

This work was supported by the Integrated Programme of SR&TD SOCA (Ref. CENTRO-01-0145-FEDER-000010), co-funded by Centro 2020 program, Portugal 2020, European Union, through the European Regional Development Fund and by the EU/EFPIA Innovative Medicines Initiative 2 Joint Undertaking under grant agreement No 806968.

## References

- Emily Chen, Kristina Lerman, and Emilio Ferrara. 2020. Tracking social media discourse about the covid-19 pandemic: Development of a public coronavirus twitter data set. *JMIR Public Health and Surveillance*, 6(2):e19273.
- Mike Conway and Daniel O’Connor. 2016. Social media, big data, and mental health: current advances and ethical implications. *Current opinion in psychology*, 9:77–82.
- Paul Ekman. 1992. An argument for basic emotions. *Cognition & emotion*, 6(3-4):169–200.
- Ethan Fast, Binbin Chen, and Michael S Bernstein. 2016. Empath: Understanding topic signals in large-scale text. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 4647–4657. ACM.
- Junling Gao, Pinpin Zheng, Yingnan Jia, Hao Chen, Yimeng Mao, Suhong Chen, Yi Wang, Hua Fu, and Junming Dai. 2020. Mental health problems and social media exposure during covid-19 outbreak. *Plos one*, 15(4):e0231924.
- Ryan CW Hall, Richard CW Hall, and Marcia J Chapman. 2008. The 1995 kikwit ebola outbreak: lessons hospitals and physicians can apply to future viral epidemics. *General hospital psychiatry*, 30(5):446–452.
- Cyrus SH Ho, Cornelia Yi Chee, and Roger Cm Ho. 2020. Mental health strategies to combat the psychological impact of covid-19 beyond paranoia and panic. *Ann Acad Med Singapore*, 49(1):1–3.
- Pedro Silva Moreira, Sonia Ferreira, Beatriz Couto, Mafalda Machado-Sousa, Marcos Fernandez, Catarina Raposo-Lima, Nuno Sousa, Maria Pico-Perez, and Pedro Morgado. 2020. Protective elements of mental health status during the covid-19 outbreak in the portuguese population. *medRxiv*.
- Norbert Müller. 2014. Infectious diseases and mental health. *Comorbidity of Mental and Physical Disorders*, page 99.
- Michael J Paul and Mark Dredze. 2017. Social monitoring for public health. *Synthesis Lectures on Information Concepts, Retrieval, and Services*, 9(5):1–183.
- Tim R Wind, Marleen Rijkeboer, Gerhard Andersson, and Heleen Riper. 2020. The covid-19 pandemic: The ‘black swan’ for mental health care and a turning point for e-health. *Internet Interventions*.
- Yu-Tao Xiang, Yuan Yang, Wen Li, Ling Zhang, Qinge Zhang, Teris Cheung, and Chee H Ng. 2020. Timely mental health care for the 2019 novel coronavirus outbreak is urgently needed. *The Lancet Psychiatry*, 7(3):228–229.
- Hao Yao, Jian-Hua Chen, and Yi-Feng Xu. 2020. Patients with mental health disorders in the covid-19 epidemic. *The Lancet Psychiatry*, 7(4):e21.