



# Multilingual topic modeling for tracking COVID-19 trends based on Facebook data analysis

Amina Amara<sup>1</sup> · Mohamed Ali Hadj Taieb<sup>2</sup> · Mohamed Ben Aouicha<sup>2</sup>

Accepted: 21 October 2020 / Published online: 13 February 2021  
© Springer Science+Business Media, LLC, part of Springer Nature 2021

## Abstract

Social data has shown important role in tracking, monitoring and risk management of disasters. Indeed, several works focused on the benefits of social data analysis for the healthcare practices and curing domain. Similarly, these data are exploited now for tracking the COVID-19 pandemic but the majority of works exploited Twitter as source. In this paper, we choose to exploit Facebook, rarely used, for tracking the evolution of COVID-19 related trends. In fact, a multilingual dataset covering 7 languages (English (EN), Arabic (AR), Spanish (ES), Italian (IT), German (DE), French (FR) and Japanese (JP)) is extracted from Facebook public posts. The proposal is an analytics process including a data gathering step, pre-processing, LDA-based topic modeling and presentation module using graph structure. Data analysing covers the duration spanned from January 1st, 2020 to May 15, 2020 divided on three periods in cumulative way: first period January-February, second period March-April and the last one to 15 May. The results showed that the extracted topics correspond to the chronological development of what has been circulated around the pandemic and the measures that have been taken according to the various languages under discussion representing several countries.

**Keywords** Social media analysis · Covid-19 · Topic modeling · Facebook · Data visualization · Multilingual

## 1 Introduction

People around the world need continuously public safety and emergency management services. These services require tools that detect quickly the occurrence of emergencies and create a correct and detailed idea about the current situation. Such tools may help alleviating desolation under harsh conditions related to natural or human-made disasters by fast and semi-automatic identification of the type, extent, place, intensity, and implications of the disaster and the knowledge transfer about the disaster issue. Content provided by the user is growing during disasters in comparison with the normal situation. Therefore, the analysis of the

provided content is necessary to identify emergencies and recent disasters, based on social networks and media search, and direct relief proportional to the needs.

Disasters are ranging from earthquakes, floods, hurricanes, droughts, tsunamis, landslides, terrorism and wars to solar flares, cosmic explosions and meteorites. These risks imperil people, populations, civilization, and humankind. Defending against these threats requires various kinds of endeavors supported by varied tools and a great range of technical and human capabilities. Advanced knowledge on the nature of emergencies and effective awareness may help to reduce the costs of the disasters. Information and communication tools are vital in modeling emergencies and population response, and in the accurate monitoring of disasters. Advanced technologies can benefit from the development of the new means and methods of data and information transmission, including the Internet, social networks and social media [1]. Various types of research analyze the benefit of social media in monitoring the defending reaction against the disasters and for making accurate strategies in real-time by responding to the urgent population needs [2–4]. They focused on the use of social media in relation with different kinds of disasters such as the earthquake [5–7], tsunami [8], typhoons [9], storms [10], flood [11] and the health risks

---

This article belongs to the Topical Collection: *Artificial Intelligence Applications for COVID-19, Detection, Control, Prediction, and Diagnosis*

✉ Mohamed Ali Hadj Taieb  
mohamedali.hadjtaieb@gmail.com

<sup>1</sup> Multimedia, Information systems and Advanced Computing Laboratory, University of Sfax, Sfax, Tunisia

<sup>2</sup> Faculty of Sciences, University of Sfax, Sfax, Tunisia

such as epidemic virus spreading [12], prediction of the contagious population behaviour and accurate detection and identification of professionally unreported drug side effects using widely available public data (open data) [13, 14].

Due to its open policies in front of information extraction, Twitter data has been exploited by the majority of works despite the fact that Facebook is ranked as the first social network having active users. The Digital 2019 reports<sup>1</sup> include extensive insights into people's use of the world's top social platforms in more than 230 countries and territories around the world. Worldwide social media users' numbers have grown to almost 3.5 billion at the start of 2019, with 288 million new users in the past 12 months pushing the global penetration figure to 45 percent. Facebook maintains its top platform ranking in early 2019. The number of monthly active Facebook users grew steadily across the past 12 months, and the platform's latest earnings announcement reports year-on-year user growth of almost 10 percent.

To our knowledge, the present work is the first that extracts insights from Facebook provided centric COVID-19 data. In fact, the majority of works on the pandemic COVID-19 are based on Twitter data. This paper describes the process of building the provided crowdsourcing trends and its evolution within the time since the first of January 2020. This study exploits the LDA-based topic modeling method in a multilingual framework and provides a novel representation method based on graph structure and handling it with graph visualization software.

The rest of the paper is organized as follows. Section 2 provides an overview of the research works focusing on studying the utility of social data in relation to different kinds of disasters by highlighting those on the health domain. Section 3 depicts some research works focusing on COVID-19 related social data analysis. The used topic modeling method, i.e. Latent Dirichlet Allocation (LDA), is detailed in Section 4. Then, the proposed Facebook data-based tracking system of the trends through time is depicted in Section 5 with its components. The multilingual gathered COVID-19 centric dataset and its statistics as well as characteristics are presented in Section 6. Section 7 reports on the interpretation of the results (COVID-19 trends) in a multilingual framework. The final section is devoted to presenting our conclusions and future research.

## 2 Related work

This section focuses on research works that have studied the use of social data on monitoring and tracking disasters in two dimensions, natural and health disasters.

<sup>1</sup><https://wearesocial.com/global-digital-report-2019>

### 2.1 Natural disasters

Social media such as Facebook and Twitter have proven to be a useful resource to understand public opinion towards real-world events. After the great east Japan earthquake in 2011, numerous tweets were exchanged on Twitter. Several studies have already pointed out that micro-blogging systems have shown potential advantages in emergency situations, but it remains unclear how people use them. [5] investigated over 1.5 million Twitter messages (tweets) for the period ranging from 9 March 2011 to 3 May 2011 in order to track awareness and anxiety levels in the Tokyo metropolitan district to the 2011 Tohoku Earthquake and subsequent tsunami and nuclear emergencies. [6] gathered tweets immediately after the earthquake and analyzed various factors, including locations. The results showed that the people in the disaster area tend to directly communicate with each other (reply-based tweet). On the other hand, people in the other area prefer to spread information from the disaster area by using retweet.

An important characteristic of Twitter is its real-time nature. For example, when an earthquake occurs, people make many Twitter posts (tweets) related to the earthquake, which enables detection of earthquake occurrence promptly, simply by observing the tweets. [7] investigated the real-time interaction of events such as earthquakes, in Twitter, and proposed an algorithm to monitor tweets and to detect a target event. To detect such an event, they proposed a classifier of tweets based on features such as the set of keywords used in a tweet, the number of words, and their context. They consider each Twitter user as a sensor. As an application, they construct an earthquake reporting system in Japan.

PEARY et al. [8] show that during the 2011 east Japan earthquake and tsunami, social media such as Twitter and Facebook served as a lifeline for directly affected individuals and a means of information sharing. Social media was used to perform vital relief functions such as safety identification, displaced-persons locating, damage information provision, support for disabled individuals, volunteer organization, fund-raising, and moral support systems. Their study discusses the potential for the public, civil society, and government organizations to utilize social media in disaster preparedness and response.

Daga [9] implemented content analysis and Social Network Analysis (SNA) on the tweets regarding different situations to be able to understand and depict a visual representation of interaction between users. They studied the user interaction of the Filipino community between two major typhoons that hit the Philippines. Results revealed that users tend to seek and share information from reliable sources such as news websites and verified Twitter users. Determining the interaction

of Twitter users in an online community plays a vital role in information dissemination and allows appropriate response during disaster and emergency situations. In this study, SNA was used to help better understand and reveal the social interaction related to typhoon Haiyan<sup>2</sup> and Hagupit<sup>3</sup>. Results of this study have shown that user interaction among Filipino online community both Haiyan and Hagupit was influenced by Twitter-verified users having the highest value betweenness centrality score.

Moreover, social media use increases in natural disasters such as cyclones, hurricanes, or typhoons. Ulvi et al. [10] investigated the roles of social media and mainstream media on hurricanes and how they may potentially have a bigger impact. They studied the influences and risk factors of media and their role in the distribution of information were observed. They concluded that social media platforms helped spread awareness, support, and warnings. Social media has shown to have impactful effects during tropical storms around the world. Public health professionals and emergency response team should utilize social media in relief for victims.

Harnessing the crowdsourced information, under social media platforms, has become an opportunity for authorities to obtain enhanced situation awareness data for efficient disaster management practices. Nonetheless, the current disaster-related Twitter analytics methods are not versatile enough to define disaster impacts levels as interpreted by the local communities. Kankanamge et al. [11] prepared a data analysis framework, and identifying highly impacted disaster areas as perceived by the local communities. For this, the study used real-time Twitter data posted during the 2010–2011 South East Queensland Floods. The findings reveal that utilizing Twitter is a promising approach to reflect citizen interests.

## 2.2 Health domain

Social media showed a great capacity to allow public participation in content creation and circulation.

Ding et al. [15] presented their study focusing on the use of social media during the H1N1 flu epidemic in the U.S. and China. This study demonstrates that governmental structures may use social media tools for one-way dissemination of risk decisions and policies. In contrast, the general public may get around the institutional control of risk information through extra-institutional collaborative risk communication to extract truths about the emerging risks.

Achrekar et al. [16] were interested in reducing the impact of seasonal influenza epidemics and other pandemics such as the H1N1 due to its paramount importance for public health authorities. Studies have shown that effective interventions can be taken to contain the epidemics if

early detection can be made. They presented the social network enabled flu trends framework, which monitors messages posted on Twitter with a mention of flu indicators to track and predict the emergence and spread of an influenza epidemic in a population. Based on the data collected during 2009 and 2010, they find that the volume of flu related tweets is highly correlated with the number of influenza-like illness cases reported by the centers for disease control and prevention.

In the context of seasonal influenza, it can cause various complications, worsen chronic illnesses, and sometimes lead to deaths. In fact, during 2009 H1N1<sup>4</sup> flu pandemic, up to 203,000 deaths occurred worldwide. Early detection and prediction of disease outbreak is critical because it can provide more time to prepare a response and significantly reduce the impact caused by a pandemic. Lee et al. [17] presented a system that predicts future influenza activities, provides more accurate real-time assessment than before, and combines real-time big social media data streams for predictive models to generate accurate predictions. Prediction of further flu levels can represent a big leap because such predictions provide insights for public health that can be benefit for planning, resource allocation, treatments and prevention.

Sharma et al. [18] studied the pandemic of Zika virus infection. More publicized and of greater concern is the epidemic of microcephaly in Brazil, manifested by an apparent 20-fold increase in incidence from 2014–2015, believed to be caused by Zika virus infection in pregnant women. The increase in the spread of the disease has caused rapid activity surrounding it in social media. They used Facebook for the dissemination of public health information via social media. For them, it is important to spread the right information that helps the public to preventative guidelines. The use of Facebook is argued to the fact that its universal availability, outreach, and substantial influence on the information available to the public.

Pruss et al. [19] examined Twitter discussion surrounding the 2015 outbreak of Zika. Their study is based on gathered data from Twitter mentioning Zika geolocated to North and South America. Using a multilingual topic model, they automatically identified and extracted the key topics of discussion across the dataset in English, Spanish, and Portuguese. They examined the variation in Twitter activity across time and location, finding that rises in an activity tend to follow major events.

Zarrad et al. [20] exploited social media such as social networks, blogs and forums to analyze users' opinions, attitudes, and emotions about news or social events. In fact, they presented a case study about MERS<sup>5</sup> virus in Kingdom of Saudi Arabia to evaluate their approach.

<sup>2</sup>[https://fr.wikipedia.org/wiki/Typhon\\_Haiyan](https://fr.wikipedia.org/wiki/Typhon_Haiyan)

<sup>3</sup>[https://en.wikipedia.org/wiki/Typhoon\\_Hagupit\\_\(2014\)](https://en.wikipedia.org/wiki/Typhoon_Hagupit_(2014))

<sup>4</sup><https://www.cdc.gov/flu/pandemic-resources/2009-h1n1-pandemic.html>

<sup>5</sup><http://www.emro.who.int/fr/health-topics/mers-cov/mers-cov.html>

Tran and Lee [21] conducted a study of understanding and mining the spread of Ebola-related information on social media. In particular, they are based on large-scale data-driven analysis of geo-tagged social media messages to understand citizen reactions regarding Ebola, build information propagation models, and analyze spatial, temporal and social properties of Ebola-related information. Their work provides new insights into the Ebola outbreak by understanding citizen reactions and topic-based information propagation, as well as providing a support for future public health crises. Missier et al. [22], also, proposed a system for tracking Dengue epidemics using Twitter content classification and topic modeling.

During the epidemics or pandemics, the potential threat to the society is the propagation of rumors through social media. Sicilia et al. [23] proposed a rumour detection system that leverages on newly designed features, including influence potential and network characteristics measures. They tested their approach on a real dataset composed of health-related posts collected from Twitter microblog.

In the next section, we revise the research works focusing on the use of social data analysis in relation to the COVID-19 pandemic.

### 3 COVID-19 related social data analysis research works

The novel coronavirus disease, named COVID-19, was identified at the first time in Wuhan, Hubei Province, China by the end of December 2019. It has rapidly outbroken worldwide leading to a global health emergency on 30 January 2020 and it has been announced as a pandemic by the World Health Organization (WHO) on 11 March 2020<sup>6</sup>. This ongoing pandemic puts all societal levels in unprecedented situation and pushed many governments around the world to enforce different measures to contain the spread of this ongoing coronavirus. Distance learning, self-quarantines and social distancing are among the maintained measures.

These enforced unprecedented measures, especially “social distancing” the most widely used of such measures, have impacted the lifestyle of people around the globe, and bring them to the frontline social media platforms for both chatting and news. Social media websites like Facebook and Twitter are playing a central and significant role, more than ever, as adequate tools that allow people to stay connected during crises for global social discussions. As more and more social interactions turn online, an important amount of conversations around this ongoing coronavirus are continuously expanding. Researchers are

mainly using these online conversations to understand the spread of this novel coronavirus, explore its aspects as well as monitoring people’s behaviours regarding to the global health emergency and so forth. Although the COVID-19 epidemic’s appearance is relatively new, there is a rapid move in the research landscape since more than 24,000 research papers<sup>7</sup> are published online. These researches are distributed over several disciplines such as social science, medicine, public health, and so on. Some of these research papers are preprint and have not yet been peer-reviewed. Table 1 presents some works taking profit from the generated social data to perform the pandemic COVID-19 related analysis.

In the computer science discipline, at the time of our writing (in mid-May 2020) and to the best of our knowledge, main publishers like Springer, IEEE Explorer, Science Direct and ACM published a few preliminary research works [24–29] about social media’s dynamics around the context of COVID-19 as it is illustrated in Table 1. However, there is prosperity of studies [30–39] that are pre-print papers investigating the evolutionary aspects of the unfolding coronavirus disease. Collected data are used to analyze the behavioral change, track COVID-19 related misinformation and rumors’ spreading [30], measure the emotional responses and worries<sup>8</sup> about the pandemic [32, 36], and so forth. Banda et al. [32] have released a dataset<sup>9</sup> of about 4.4 million of daily tweets related to coronavirus context collected through the Twitter stream API using keywords as ‘Covid-19’ and ‘Coronavirus’. Chen et al. [30] proposed a multilingual Covid-19 Twitter dataset, which is made available to the research community<sup>10</sup>, collected since January 2020 using Twitter stream API and Tweepy to follow specific keywords and trending subjects simultaneously. This dataset helps tracking coronavirus-related unverified rumors, enabling the understanding of users’ sentiment towards this global crisis and more. Alqurashi et al. [31] proposed an Arabic COVID-19 Twitter dataset presenting preliminarily analysis on COVID-19 tweets to reveal arabic users’ behaviour and sentiment during this pandemic. Other works [34, 35] are focusing on related COVID-19 pandemic data shared on multiple social media platforms. Indeed, Gao et al. [35] have released a multilingual dataset of social media posts extracted from the two micro-blogging websites Twitter and Weibo in 3 different languages: English, Japanese and Chinese language. Last but not least, to the best of our knowledge and at

<sup>6</sup><https://www.who.int/dg/speeches/detail/who-director-general-s-opening-remarks-at-the-media-briefing-on-covid-19—11-march-2020>

<sup>7</sup><https://www.technologyreview.com/2020/03/16/905290/coronavirus-24000-research-papers-available-open-data/#Echobox=1585947735>

<sup>8</sup><https://github.com/ben-aaron188/covid19worry>

<sup>9</sup><https://zenodo.org/record/3738018#.XrALQ1UzZ0w>

<sup>10</sup><https://github.com/echen102/COVID-19-TweetIDs>

**Table 1** Source of collected data and type of language used in related works

| Collected data | Group       |              | Social media platform                       |
|----------------|-------------|--------------|---|
|                | Monolingual | Multilingual |   |
| [31]           | X           |              | Twitter                                     |
| [24]           |             | X            | Twitter                                     |
| [30]           |             | X            | Twitter                                     |
| [36]           | X           |              | Twitter                                     |
| [38]           |             | X            | Twitter                                     |
| [40]           |             | X            | Facebook                                    |
| [33]           | X           |              | Facebook                                    |
| [39]           |             | X            | Instagram                                   |
| [35]           |             | X            | Twitter And Weibo                           |
| [34]           | X           |              | YouTube, Instagram,<br>Twitter, Reddit, Gab |

the time of this writing, there are currently very few works focusing on Covid-19 Facebook data analysis [33, 37, 40]. [37] examine the relationship between the geographic spread of COVID-19 and the geo-location information from users of social networks such as Facebook, in the United States and Italy. Perrotta et al. [40] proposed a rapid response monitoring system through continuously running survey across eight countries. They collected key information on people's health status, attitudes and close social contacts by recruiting participants through targeted Facebook advertisement campaigns. Moreover, based on Facebook data, a dashboard for Good Mobility<sup>11</sup> is conceived by a group of infectious disease epidemiologists, at universities around the world, to provide daily updates to decision-makers on how people are moving and where they live, in order to help health organizations to improve the effectiveness of their campaigns about the epidemic response.

## 4 Topic modeling

Latent Dirichlet Allocation Topic modeling is one of the unsupervised machine learning methods, widely used in natural language processing, used for discovering hidden semantic structures, known as “topics”, in a text document. Topics mean the hidden relations that link words in a vocabulary and their occurrences in documents. Topic modeling [41] seeks to find key concepts throughout the whole corpus and annotate the documents of the corpus based on these concepts. It provides a useful view of a large corpus in terms of individual documents and relationships between them. Latent Dirichlet Allocation (LDA), proposed by [42], is one of the most popular and recent topic

modeling techniques. It is an unsupervised probabilistic generative technique for modeling text documents in a given text corpora as mixtures of latent topics based on Bayesian models. Each document is represented as a probabilistic distribution over latent topics and that per-document topic distributions share a common Dirichlet prior. Each topic in the LDA model is defined as a probabilistic distribution over words, and those per-topic word distributions share a common Dirichlet prior as well. Figure 1 shows the LDA topic modeling process.

Given a collection of  $M$  documents, each document  $d$  is composed of  $N_d$  words, with  $d \in \{1, \dots, M\}$ . In order to model the collection of documents, the LDA generative process [41] is executed as follow:

- For each topic  $t$  ( $t \in \{1, \dots, K\}$ ), select a multinomial distribution  $\phi_t$  whose hyper-parameter  $\beta$  follows the Dirichlet distribution.
- For each document  $d$  ( $d \in \{1, \dots, M\}$ ), select a multinomial distribution  $\theta_d$  having an hyper-parameter  $\alpha$  which follows the Dirichlet distribution.
- For each word  $w_n$  ( $n \in \{1, \dots, N_d\}$ ) in document  $d$ ,
  - i) Choose a topic  $z_n$  from  $\theta_d$
  - ii) Choose a word  $w_n$  from  $\phi_{z_n}$

All words included within  $M$  documents are observed variables while the other components, composed of topics  $\phi_t \forall t \in \{1, \dots, K\}$  the per-document topic distribution  $\theta_d \forall d \in \{1, \dots, M\}$  and the per-word topic distribution are not known. These latter items are denoted as hidden variables which are predicted from the analysis of observed variables, i.e. data. The last two variables  $\alpha$  and  $\beta$  are hyper-parameters. The corpus probability is expressed as follow:

$$p(D | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \times \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d\theta_d \quad (1)$$

Using LDA-based topic modeling, terms in the set of documents are used to generate vocabulary which is used to discover hidden topics in a large corpus. Document is represented as a mixture of topics where a topic is defined by a probability distribution over the set of terms. In our case, the corpus is composed of Facebook public posts so that each post represents a text document. LDA-based topic modeling is applied on them in order to discover the main discussed topics during COVID-19 disease.

In fact, a topic model is a kind of a probabilistic generative model. As described above, the goal of topic modeling is to automatically discover the topics in a collection of documents. The documents themselves are examined, whereas the topic structure—the topics, per-document topic distributions, and the per-document per-word topic assignments—is

<sup>11</sup><https://visualization.covid19mobility.org/>

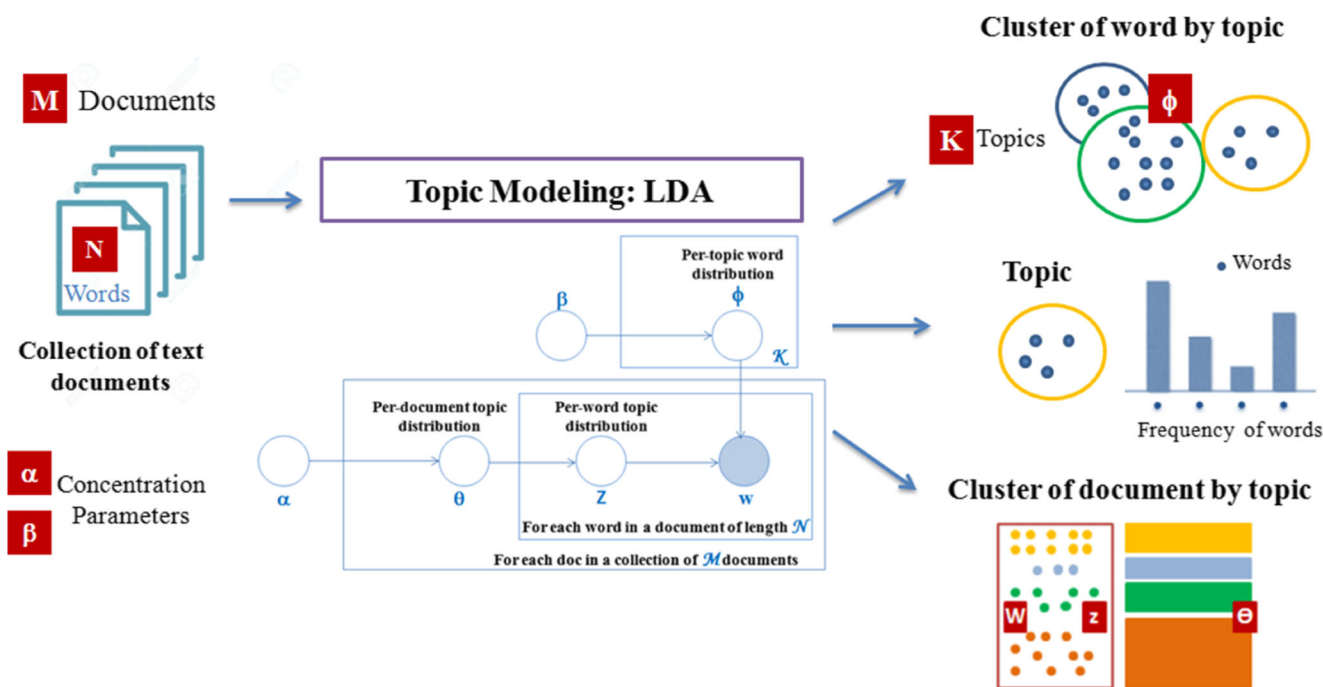


Fig. 1 LDA-based topic modeling process

hidden structure. Given a corpus, the goal is estimating the posterior distribution of unknown model parameters and hidden variables. The advantage of LDA is that the document-generative process can be adapted to other kinds of analysis, keeping only the analogy between document-topic-word and other kinds of objects. Therefore, it can be exploited for building predictive models [43–45]. [43] used natural language processing and LDA, to identify reliable patterns from within research narrative documents to distinguish studies that complete successfully, from the ones that terminate. Damevski et al. [45] exploited LDA-based topic modeling for predicting future developer behaviour in the integrated development environment. Therefore, the LDA can be used as predictive model to detect novel trends in the behaviour of the social media users’ topics or to give a deeper analysis about a specific topic by detailing the documents providing the target topic. Also, it is important to note that LDA can be considered as a fuzzy classification because it provides a soft degree of belonging of the documents to a specific topic.

### 5 Facebook-based COVID-19 tracking trends system

In this section, we describe the proposed system including its different components and tools. As it is already mentioned, the system extract multilingual information. Figure 2 illustrates the different components from the collection and posts extraction to the topic modeling and visualization.

#### 5.1 Selecting active Facebook users and extracting public posts

The starting dataset covers Facebook users from the whole world. Users are collected during our previous research work [46]. The Facebook users are examined for identifying active users based on the publication rate since first January 2020. The selected users are candidates for having published on COVID-19 context.

The candidate users are explored for scrapping the published posts about the COVID-19 by checking the existence of the built vocabulary. Then, each post is stored after extracting the URL from the post. In fact, the Facebook user can provide his own content, share a content with a comment or in its initial format. So, these posts will be used for building an incremental idea about the COVID-19 along the time period starting from the first January 2020. As it is known, the politic access to the information is so close when it is compared to Twitter as example. Therefore, we design some patterns to locate the public posts with the time related information. The target posts are selected based on multilingual COVID-19 vocabulary built using machine translation.

#### 5.2 Multilingual topic modeling

The proposed work is focused on multilingual context of COVID-19. In fact, the keyword *COVID-19* is exploited in its written format by several languages using Latin letters. The gathered data pertains to users in the whole world (see

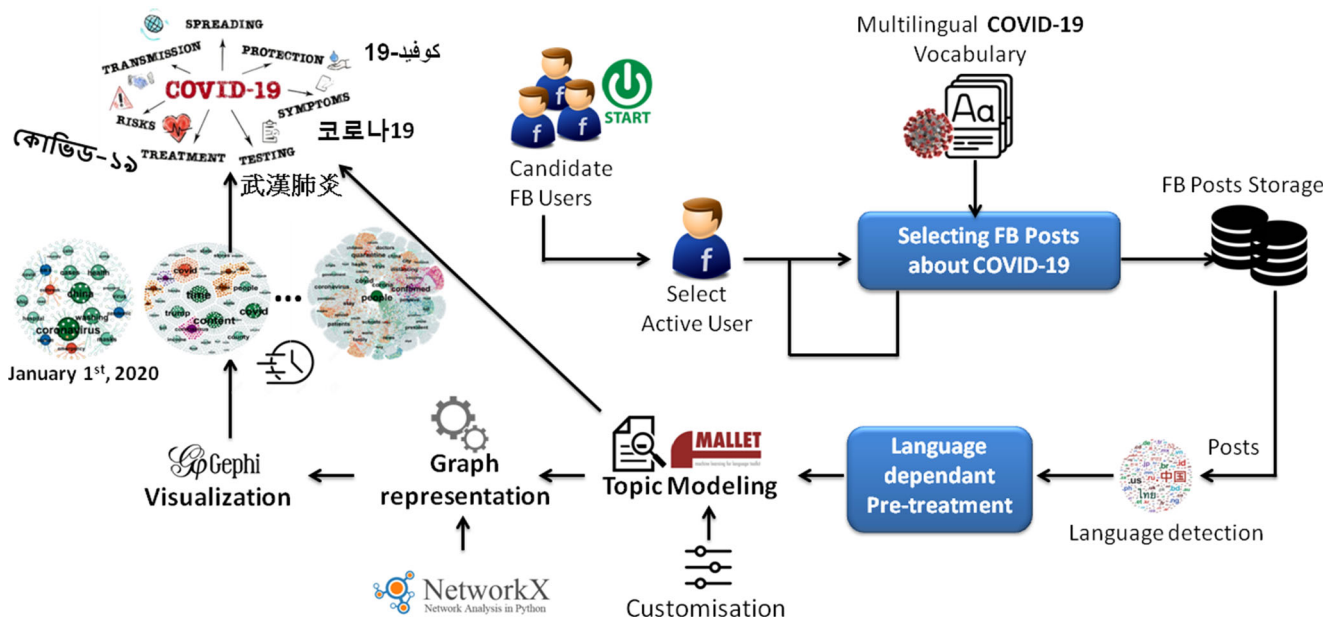


Fig. 2 Facebook-based COVID-19 tracking trends evolution system

Fig. 3) in different languages. The language and the target url pages are identified. We exploit the language detection library implemented in plain Java and covering more than 50 languages. This library<sup>12</sup> is open source Apache license 2.0 and it calculates language probabilities from features of spelling through naïve bayes with character n-gram. It generates the language profiles from training groups. A profile is the set of probabilities of all spelling in each language.

We exploit the Latent Dirichlet Allocation (LDA) [42] that it is considered as “generative probabilistic model” of a collection of composites made up of parts. Its uses include mainly topic modeling. The composites are documents and the parts are words and/or phrases. In our study, the documents are the published posts.

The computing process of the probabilistic topic model estimated by LDA consists of two tables (matrices). The first table describes the probability of selecting a particular word when sampling a particular topic. The second table describes the chance of selecting a particular topic when sampling a particular post. The LDA algorithm is composed of the following detailed steps in relation with our context:

- i. Select the unique set of words.
- ii. Select the set of posts according to a specific language.
- iii. Fix how many parts you want per posts (sample from a Poisson distribution).

- iv. Specify the number of topics as outputs.
- v. Affect a number between not-zero and positive infinity to the parameter alpha.
- vi. Affect a number between not-zero and positive infinity to the parameter beta.
- vii. Build the ‘words-versus-topics’ table. For each column, the beta is used as input for the Dirichlet distribution (which is a distribution of distributions). Each sample will fill the columns of the table and give the probability of each word per topic (column).
- viii. Build the ‘posts-versus-topics’ table. For each row, the parameter alpha is used by Dirichlet distribution as the input. Each sample will fill the rows of the table and give the probability of each topic (column) per posts.
- ix. Build the actual posts. For each of them, a) look up its row in the ‘posts-versus-topics’ table, b) sample a topic based on the probabilities in the row, c) go to the ‘words-versus-topics’ table, d) look up the topic sampled, e) sample a part based on the probabilities in the column, f) repeat from step 2 until reaching how many words the post was set to have.
- x. Moreover, the number of words representing the topic can be fixed at the beginning of the algorithm. The number of topics can be viewed as a number of clusters and the probabilities as the pertaining degree to the cluster. LDA process plays the role of soft-clustering between posts and the words composing the topics.

Posts follow a pre-processing step for preparing data to the topic modeling process. In fact, this includes mainly the remove of html tags and the stop words according to

<sup>12</sup><https://code.google.com/archive/p/language-detection>



**Fig. 3** Distribution of the COVID-19 related data through countries around the world

its own language. Mallet<sup>13</sup> is the tool exploited in this study to realise the topic modeling task [47]. MALLET is a Java-based package for statistical natural language processing, document classification, clustering, topic modeling, information extraction, and other machine learning applications to text. Many of the algorithms in MALLET depend on numerical optimization. In fact, unsupervised topic modeling is useful for analyzing large collections of unlabeled text. The MALLET topic modeling toolkit contains sampling-based implementations of Latent Dirichlet Allocation. The toolkit is open source software, and is released under the common public license.

The topic modeling module provides as output the topic composition of posts  $P_{Topics}$  and the top  $k$  words for each topic among  $N$  topics (where  $N$  and  $k$  are predefined)  $T_{words}$ . This output can be useful for checking that the model is working as well as displaying results of the model. In addition, it reports the Dirichlet parameter of each topic. If hyperparameter optimization is turned on, this number will be roughly proportional to the overall portion of the collection assigned to a given topic.

### 5.3 Building graph representation

The two matrices provided by Mallet are used for providing a graph representation allowing us in next step to handle the topics-words based graph through graph processing tools.

<sup>13</sup><http://mallet.cs.umass.edu/index.php>

This module is implemented using python and NetworkX<sup>14</sup> which is a Python package for the creation, manipulation, and study of the structure, dynamics, and functions of complex networks and it is distributed as open source 3-clause BSD license. This module generates a file GEXF<sup>15</sup> (Graph Exchange XML Format). This graphic data can be exchanged from one application to another due to *.gexf*. Various operating systems can be used to open *.gexf* files and it only requires appropriate applications such as Gephi and gexf4j. This step provides two types of graph representation to the topic modeling output. For the former, it provides only the connection between the post and its corresponding topic based on the highest contribution of the composite to the topic. For the second, an edge is built between the topic and each post according to the participation weight to form the topic. We choose the graphic representation of the topic modeling output to express in meaningful way the evolution of the COVID-19 centric interests since the beginning of the 2020 year and in multilingual framework. The GEXF file is treated by the Gephi<sup>16</sup> as open graph viz platform for offering good specialization layout and a range of specific algorithms for graph handling and topological parameters extraction. The topics are represented by one word by customizing

<sup>14</sup><https://networkx.github.io/>

<sup>15</sup>GEXF is a language for describing complex network structures, their associated data and dynamics. This extension is used to describe files containing graphic and visualization data.

<sup>16</sup><https://gephi.org/>



the Mallet tool to provide topics with 3 words. Next, to resolve the problem of exploiting the same word in different contexts due to its polysemous nature, we visit the topics

from the higher probable to the last one, and we explore the words used to represent the topic until finding not used word by the already visited topics.

---

**Algorithm 1** Generate Graph GEXF Topic Modeling.

---

**Input:**  $P_{Topics}$  the contributions of Facebook Posts through the extracted topics,  $T_{words}$  where each topic among  $N$  extracted topics is represented by  $K$  words

**Output:**  $GraphPostsTopics$  the graph structure in format GEXF

```

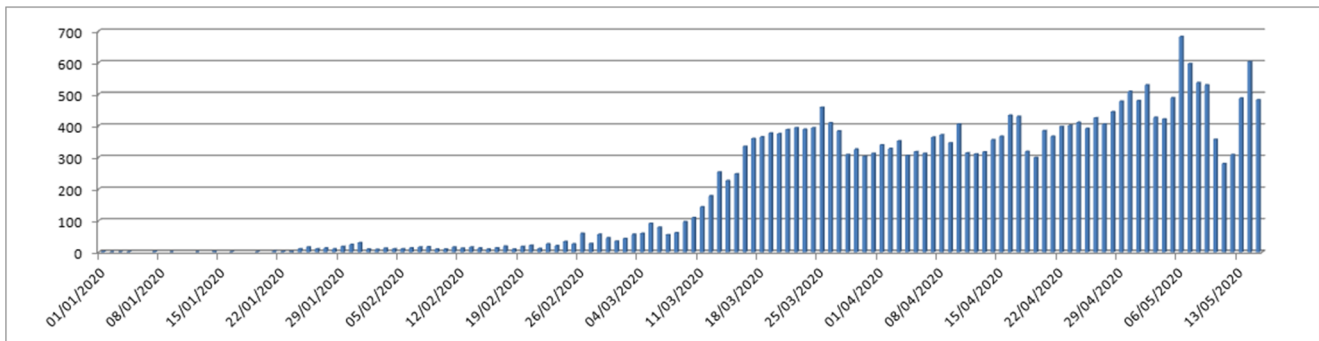
1:  $T1_{words} \leftarrow selectSingleWordTopic(T_{words})$ 
2:  $weights \leftarrow read(P_{Topics})$ 
3:  $topics \leftarrow read(T1_{words})$ 
4:  $post\_topic\_weights \leftarrow Dictionary()$ 
5: for  $weight \in weights$  do
6:    $post\_topic\_weights \leftarrow calculate\_edge\_weights(post\_topic\_weights, weight[0], weight[1], weight[2])$ 
   //weight[0]: id of the Post
   //weight[1]: path of the Post file after its pre-processing
   //weight[2]: series of values representing the participation of the Post to each topic
7: end for
8:  $GraphPostsTopics \leftarrow networkX.createGraph()$ 
9: for  $post \in post\_topic\_weights.keys()$  do
10:    $GraphPostsTopics.add\_node(post)$ 
11: end for
12: for  $topic \in topics$  do
13:    $GraphPostsTopics.add\_node(id=topic[0], label=topic[2], vizualisation=topic[1])$ 
   //topic[0]: number of the topic among  $K$  topics
   //topic[1]: topic weight
   //topic[2]: word representing the topic
14: end for
15: for  $post \in post\_topic\_weights.keys()$  do
16:    $idTopicMax \leftarrow 0$   $\triangleright$  it contains the  $id$  of the topic which  $post$  having the highest participation in it
17:   for  $tid \in post\_topic\_weights[post].keys()$  do
18:     if  $post\_topic\_weights[post][tid] > post\_topic\_weights[post][idTopicMax]$  then
19:        $idTopicMax \leftarrow tid$ 
20:     end if
21:      $GraphPostsTopics.add\_edge(idPostMax, post, weight = post\_topic\_weights[post][idTopicMax])$ 
22:   end for
23: end for

```

---

Algorithm 1 depicts the different steps for building the graph GEXF using the output of the multilingual topic modeling module. In fact, Algorithm 1 receives as inputs the  $P_{Topics}$  representing the distribution of Facebook Posts on the extracted topics and  $T_{words}$  containing the  $N$  extracted topics and the  $K$  words for each topic. The function *selectSingleWordTopic* at Instruction 1 selects the single word representing the topic among  $K$  words. The process is based on ranking the topics according to their weights. Then, it selects the word does not exist among the already selected word by traversing topics from the highest weighted topic to the last one. The entities *weights*

(Instruction 2) and *topics* (Instruction 3) are triplets as described within the algorithm. The dictionary (Instruction 4) is a general-purpose data structure for storing a group of objects. A dictionary has a set of keys and each key has a single associated value. When presented with a key, the dictionary will return the associated value. The dictionary *post\_topic\_weights* is filled using the for loop (instruction 5) such as the Keys are the ids (*weight[0]*) of the posts and Values are the weights contributions of a post to the extracted topics. The building of the graph *GraphPostsTopics* is based on adding nodes (Instructions 10 and 13) and linking them with edges (Instruction 21).



**Fig. 4** Distribution of the gathered Facebook public posts through the time since January 1st, 2020 to May 15, 2020

## 6 Dataset presentation

This section is devoted to present the multilingual dataset of posts extracted from Facebook. This dataset is in continuous evolving and can be followed through the github web site<sup>17</sup>. This presentation covers several dimensions which are detailed in next paragraphs.

### 6.1 Dataset statistics

The gathered data covers the period since January 1st, 2020 to May 15th, 2020. Figure 4 shows the distribution of the number of posts in relation with their date of posting. The first two months characterize the beginning confrontation with the new virus started from Wuhan. The virus, back then, is not known and the disease is not spread around the world. In fact, it is just localised in China. Therefore, the gathered data is few except in Japan because the disease arrived early. After the first stage, there have been an explosion in the number of posts coordinately with the appearance of the coronavirus crisis in Italy.

Figure 3 presents the geo-location coverage of the users. We tried to cover the maximum of countries where people practice different languages.

### 6.2 Multilingual aspect

Table 2 shows the number of gathered posts in relation with each language in a cumulative way. As shown in Fig. 3, the languages do not mean a specific country, but the most publications may belong to a particular country such as English from United States, French from France, Italian from Italia, Spanish from Spain and Dutch from Germany. As for Arabic is from several countries. It is important to note that the same users who provided few posts about COVID-19 (except for English and Japanese) during January and February, are the same who provide this

respectful number of posts on the next period. But then the number of posts increased significantly as the crisis worsened, and this depicts the degree of anxiety that people have reached.

The multilingual topic modeling is based on specific lists of stop words for each language: English<sup>18</sup>, Arabic<sup>19</sup>, Spanish<sup>20</sup>, Italian<sup>21</sup>, German<sup>22</sup>, French<sup>23</sup> and Japanese<sup>24</sup>.

### 6.3 Facebook users COVID-19 related behaviour

Figure 5 shows the percentages of the number of posts on COVID-19 in relation to the percentage of Facebook users who were the subject of the research who achieved a percentage equal to or greater than the abscissa  $x$ . In fact, 16.06% of users have 50% of their public posts about COVID-19. We have recorded the highest number 541 of posts on COVID-19 among 1245 posts during the studied period.

On another point, and a return to the users whom we singled out for the analysis, we find that 2% of them returned to using their Facebook accounts after their last public post before January 1st, 2020 dated back to the period between 2012 and 2018. This indicates the state of alert and suspicion of the pandemic, which prompted a number of users to return to follow and to share the news about the virus on social networking pages, especially with quarantine.

The pearson correlation between the number of COVID-19 centred posts and the total number of published posts

<sup>18</sup><https://gist.github.com/sebleier/554280>

<sup>19</sup><https://github.com/mohataher/arabic-stop-words/blob/master/list.txt>

<sup>20</sup><https://github.com/stopwords-iso/stopwords-es/blob/master/stopwords-es.txt>

<sup>21</sup><https://github.com/stopwords-iso/stopwords-it>

<sup>22</sup><https://github.com/stopwords-iso/stopwords-de/blob/master/stopwords-de.txt>

<sup>23</sup><https://github.com/stopwords-iso/stopwords-fr/blob/master/stopwords-fr.txt>

<sup>24</sup><https://github.com/stopwords-iso/stopwords-ja/blob/master/stopwords-ja.txt>

<sup>17</sup><https://mohamedalihadhtaieb.github.io/Covid19-based-Facebook-Research/>

**Table 2** Distribution of the collected COVID-19 related Facebook posts through the time and in relation with treated languages

|           | January - February 2020 | March 2020 | April 2020 | 15 May 2020 |
|-----------|-------------------------|------------|------------|-------------|
| <b>EN</b> | 261                     | 3870       | 8541       | 10294       |
| <b>AR</b> | 55                      | 817        | 2124       | 2907        |
| <b>DE</b> | 7                       | 288        | 819        | 1546        |
| <b>FR</b> | 12                      | 271        | 1021       | 2267        |
| <b>ES</b> | 14                      | 338        | 1012       | 1594        |
| <b>IT</b> | 44                      | 504        | 1394       | 1981        |
| <b>JA</b> | 230                     | 707        | 1383       | 1790        |

since January 1st, 2020 is equal to 0.62 which leads to significant correlation between the two parameters. This means that for the majority of users the number of COVID-19 posts follows the increasing of the number of total posts.

Missier et al. [48] showed that social media analytics can be used to pinpoint individuals who are actively contributing to social discourse on the specific topic of the Zika virus and its consequences, and are thus likely to be sensitive to health promotion campaigns, by focusing on Twitter content related to the Zika virus and its effect on people. For our case, also, we find some users who can be qualified as COVID-19 publish engine despite that their accounts are for persons. In fact, 1.40% of users, considered in this research, published more than 200 public posts on COVID-19 since January 1st, 2020. They participate with 24.46% of the gathered data.

## 7 Analysis of COVID-19 trends

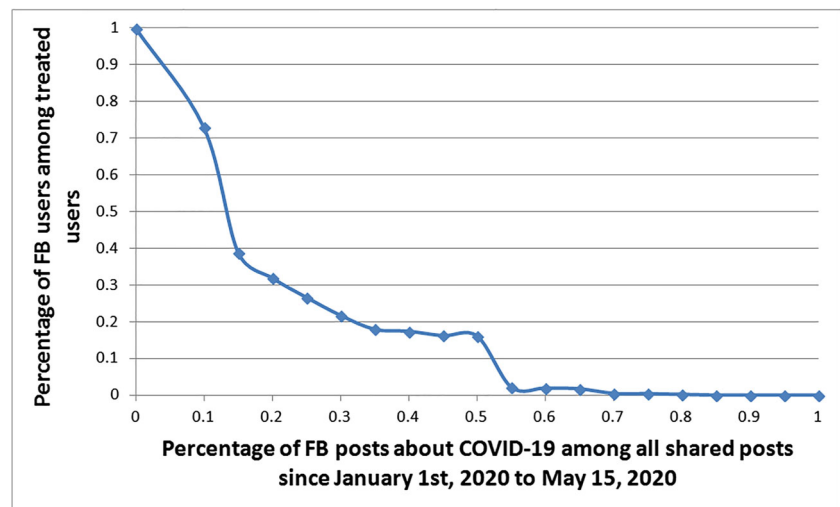
We present in this section the results of the proposed architecture for tracking the crowdsourcing COVID-19 centred trends based on social data provided by Facebook users about this pandemic. The discussion is focused on

7 different languages (English (EN), Arabic (AR), French (FR), Spanish (ES), Italian (IT), German (DE) and Japanese (JA)) and it explores the evolution over time of the users' interests from the start (January 1, 2020) of the coronavirus crisis until May 15, 2020. The differences figured in the countries represented by the languages will be highlighted. Moreover, the interests differ along with the diversity of stakeholders during this serious crisis having launched an emergency status. The results are visualized in graph structure according to the importance degree of the extracted topics linked to the posts having mostly participated to form a specified topic. The study is divided mainly on 3 periods according to the outbreak of the virus SARS-COV-2 around the world. The first period covers January and February and it is characterized by the acquaintance stage, the second period along the months March and April which includes the shock stage and radical change in daily life, and for the third until 15 May 2020 with the beginning of returning to normal life taking into consideration the coexistence with the virus.

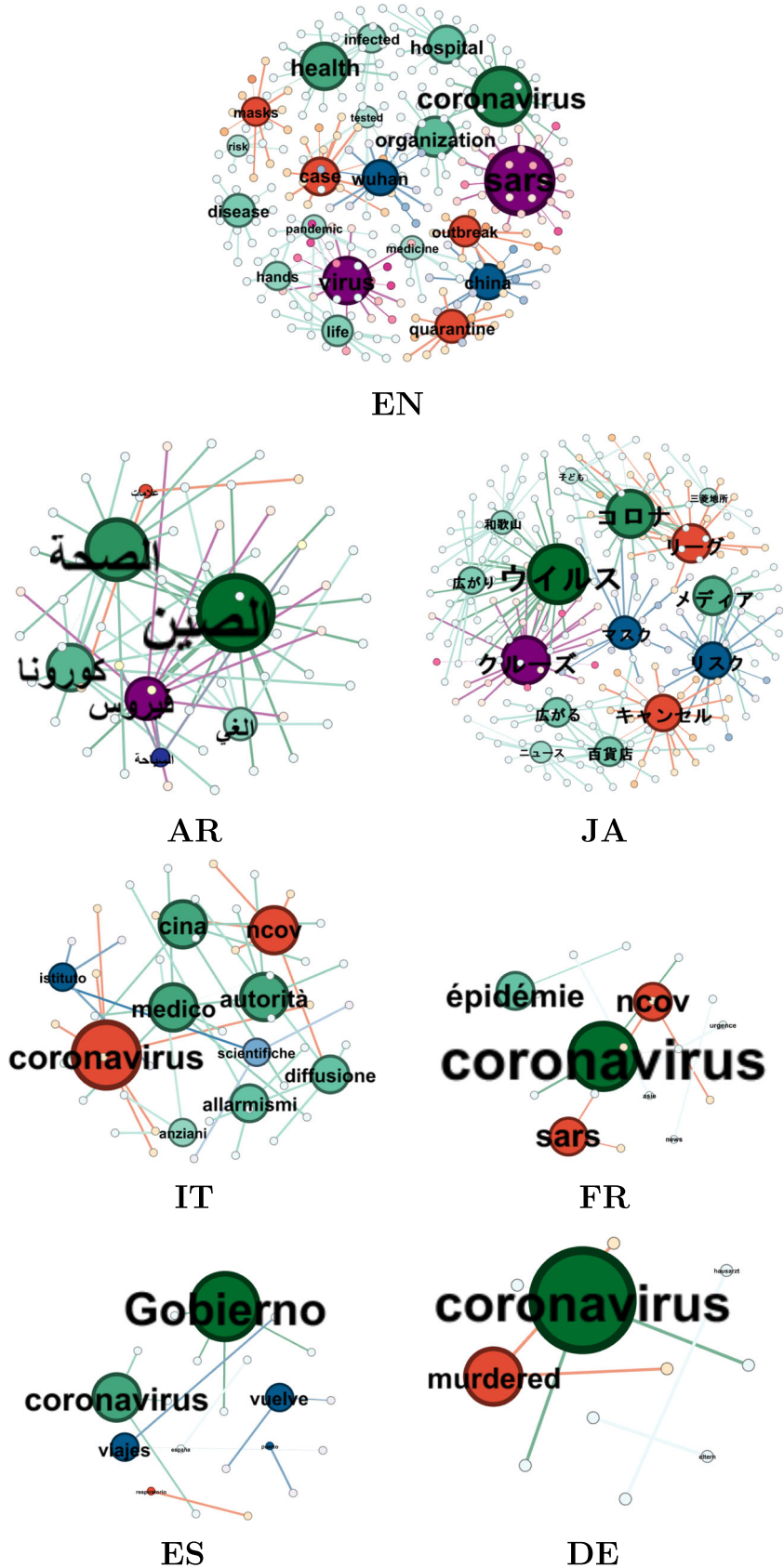
### 7.1 Analysis of first period: January-February 2020

Figure 6 shows the first topics involved in the social network Facebook in relation with different languages.

**Fig. 5** Correlation between percentage of COVID-19 Facebook posts during the period from January 1st, 2020 to May 15, 2020, and the percentage of users having post more than the fixed percentage



**Fig. 6** Facebook-based crowd-sourced COVID-19 trends covering 7 languages (EN, AR, FR, ES, IT, DE and JA) for the period January and February 2020



Mainly, users start exploiting the virus name [*Coronavirus* (EN), *كورونا* (AR) and *コロナ* (JA), *sars* and *ncov*]. A future emergency status is detected which would affect all countries of the world [*allarmismi* (IT), *urgence* (FR)]. Everyone was afraid of the deadly disease, and this was evident in the results regarding Italy, which has struck the world with the number of dead people before the disease spread to other countries [*murdered* (DE)]. Moreover, they follow the mobilization of the official authorities [*Gobierno* (ES), *autorità* (IT)] especially the health sector [*الصحة* (AR), *medico* (IT), *hausarzt* (DE)]. They talk also about its origin China and especially Wuhan [*الصين* (AR), *asie* (FR), *cina* (IT)]. Users are also watching its outbreak [*épidémie* (FR)]. Facebook surfers and parents [*“eltern”* (DE)] were observing the movement of travel for fear of the spread of the disease and the closure of land and sea borders as a result of the pandemic of coronavirus [*السباحة* (AR), *viajes, vuelve, puerto* (ES), *クルーズ* (JA)]. Infection has spread between countries through the movement of travelers. This created a global crisis and the desire of the various countries to close their borders and led to the emergence of the problem of the stuck with the lack of freedom of movement. Authorities around the world had resorted to canceling [*キャンセル* (JA), *الغبي* (AR)] many events. Eastern Asia countries (Japan and South Korea) are the first places where the infection started early which gives explanations to the quite richness of information for the English and Japanese languages. Therefore, they discuss the means of protection [*マスク* (JA), *masks, hands* (EN)] with the medical procedures [*quarantine* (EN), *キャンセル* (JA)] that were followed and began to talk about social separation.

## 7.2 Analysis of second period: March-April 2020

Regarding the second stage of the pandemic, after the ability of East Asian countries such as Japan and South Korea to control the spread of the disease, the world was shocked by the number of victims in Italy (and in a future stage in the United States of America) [*イタリア* (JA)] with the beginning of the rapid spread in France and Spain and the appearance of the first cases in most Arabic countries. All of this led to health measures taken by countries through quarantine, the suspension of air traffic and most economic activities, with the cancellation or postponement of most events.

The second stage was marked by a state of alert with the official authorities, with a general feeling of the seriousness of the situation and the exacerbation of the crisis [*リスク* (JA), *emergenza* (IT), *crisis* (ES), *crise* (FR), German *ausnahmestandard, coronakrise, gefährlicher* (DE)]. The

most shocked to people is the ability of the virus COVID 19 to inflict casualties [*Italianemergency* (IT), Spanish *fallecido* (ES), *وفاة* (AR), *morts, décès, deuil* (FR), *deaths* (EN)], and therefore we find that this topic was addressed through Facebook browsers with the names of some of the well-known persons on a national or global level. Figure 7 depicts the chronological development of subjects and the trends that people transmit about the virus in the various languages that represent many countries. It is not only interested in the months of March and April, but focuses on the cumulative nature of the discourse extending from the beginning of the year 2020 to the mentioned time limit.

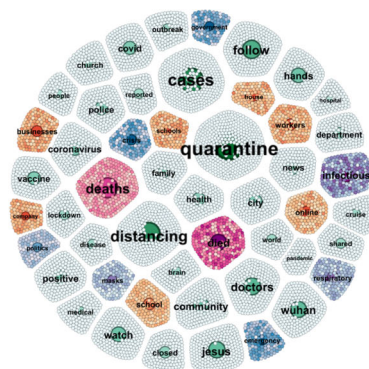
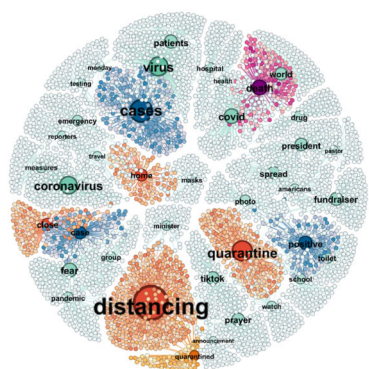
It also appears in our analysis, which is consistent with the reality that we lived through, that countries mobilized their political and financial capabilities to confront the pandemic, and societies felt that they were in a state of war, and this stimulated the use of vocabulary from the glossary of battles, as was shown in various languages [*مجاهدة, ملافحة* (AR), *に対し* (JA), *lucha, ejército* (ES), *lutte* (FR)]. Citizens also pursued their governments, politicians, and leaders, and this justifies the emergence of issues that concern governments and the political class [*大統領* (JA), *sindaco* (IT), *presidente, gobierno* (ES), *ministre, gouvernement, autorités* (FR), *behörde, bund, bundesregierung, regierung* (DE), *الحكومة, وزير, وزارة* (AR), *president, government, politics* (EN)].

One of the most important measures taken to combat the virus is quarantine [*confinement* (FR), *العزل* (AR), *quarantäne* (DE), *cuarentena, confinamiento* (ES), *quarantena* (IT), *lockdown, quarantine* (EN)], which requires people to stay in their homes to limit the spread of the disease. In fact, restricting mobility is a primary method being used to slow down the spread of COVID-19. Therefore, the topics that focus on this axis are considered one of the most exciting topics in all the languages studied in this research [*casa* (IT), *casa* (ES), *دارك, البيت* (AR), *house* (EN)].

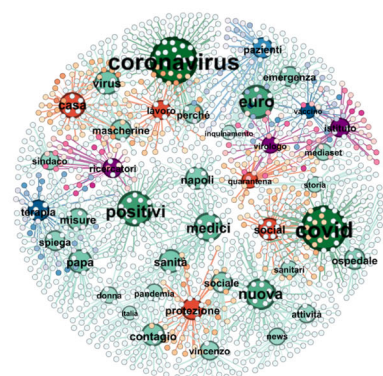
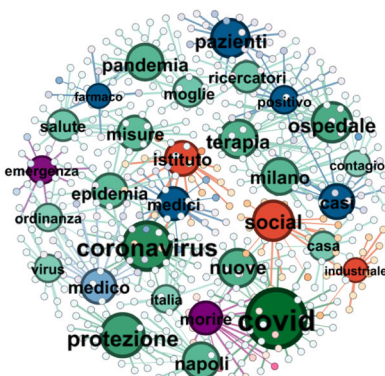
It is also worth noting children [*العبال* (AR), *子ども* (JA), *niños* (ES)] as a common topic between different browsers in different languages. This is mainly related to the quarantine procedure, which resulted in the suspension of studies for all educational levels and making parents busy with the end of the school year [*التربية, جامعة* (AR), *小学校* (JA), *estudio* (ES), *écoles, étude* (FR), *school* (EN)]. Hundreds of countries have implemented nationwide school closures and many other countries are implementing localized school closures. UNESCO<sup>25</sup> estimates that these

<sup>25</sup><https://en.unesco.org/covid19/educationresponse>

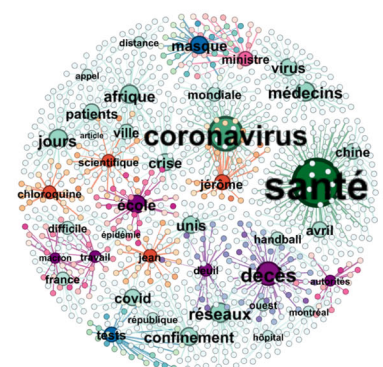
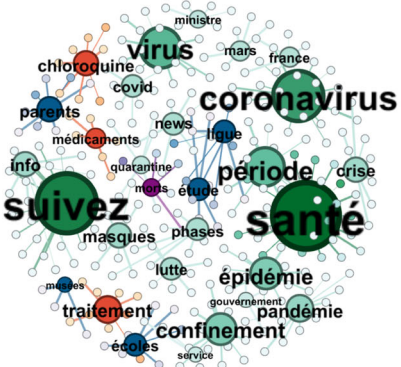
**Fig. 7** Facebook-based crowd-sourced about COVID-19 trends covering 7 languages (EN, AR, FR, ES, IT, DE and JA) for the period March (left part) and April 2020 (right part)



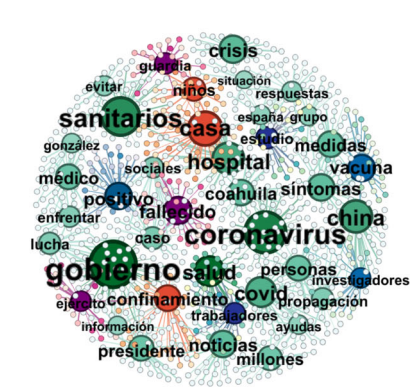
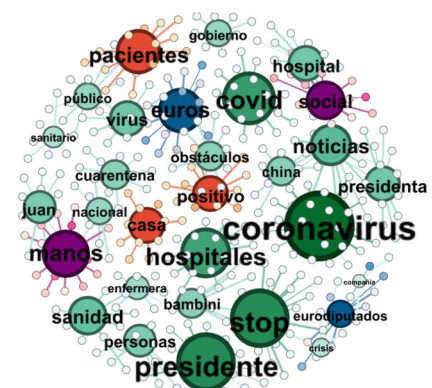
EN



IT



FR



ES

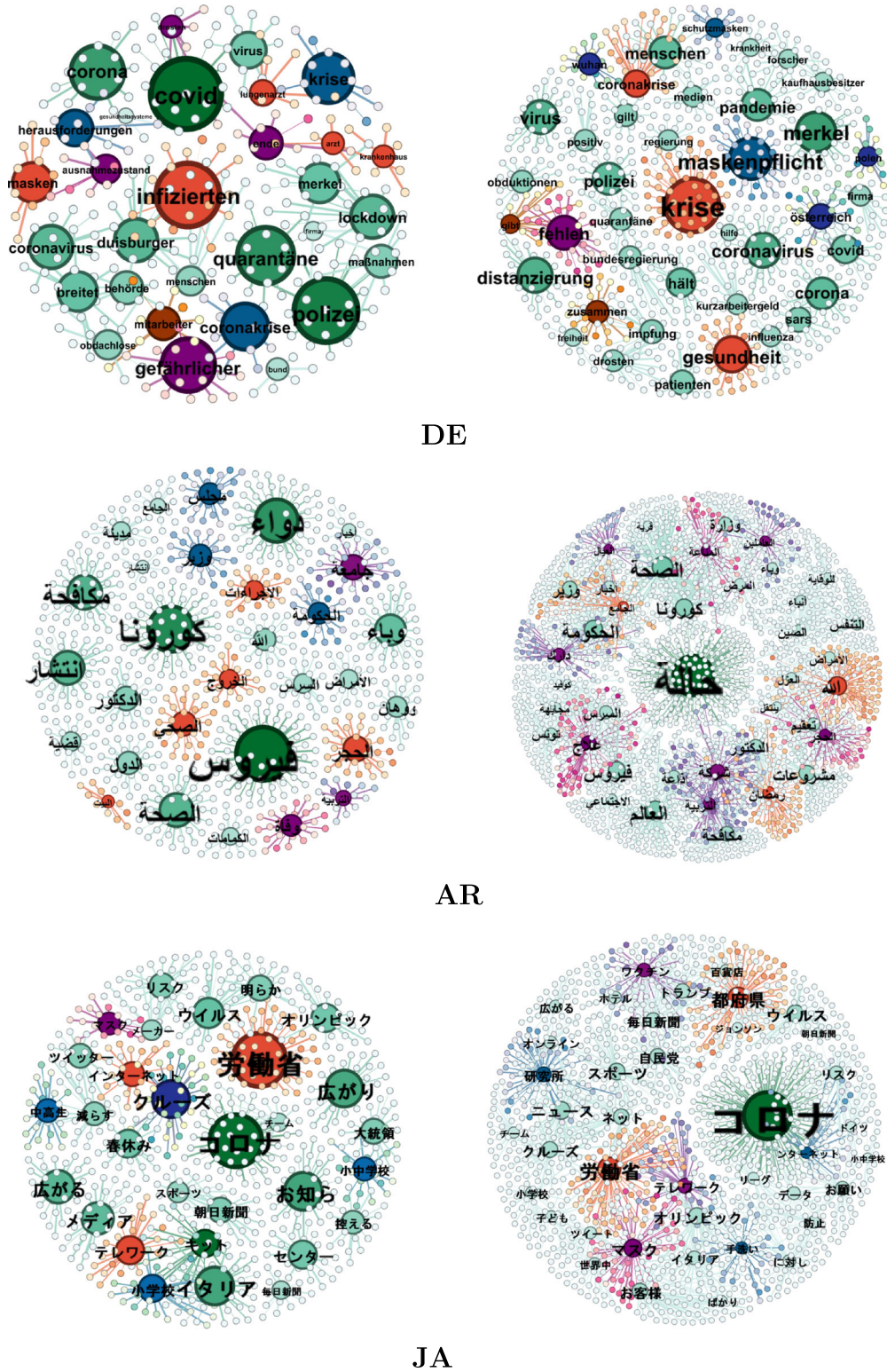


Fig. 7 (continued)

closures mean approximately 1.2 billion students, roughly 74% of all enrolled students worldwide, are experiencing a disruption in their education.

In addition to what has been mentioned, we notice the presence of topics in relation to what may be called a stock exchange of cases [*casi* (IT), *caso* (ES), *حالة* (AR), *tests*

(FR)] mainly related to medical analyzes (called tests) that are conducted and whose result is positive due to their pregnancy or disease negative [*positiv* (DE), *positive* (EN), *positivo* (ES), *positivi* (IT)]. These data were linked to several systems that were developed to monitor the spread of the disease around the world. It is noticed that the users of social media platforms, especially the Facebook in question, tend to follow the number of patients and victims and they share this information and comment on it.

In the context of searching for appropriate drug [*médicaments, traitement* (FR), *drug* (EN), *farmaco, terapia* (IT), *دواء* (AR)] and vaccination *ワクチン* (JA), *vaccino* (IT), *invacuna* (ES), *impfung* (DE), *vaccine* (EN)] to prevent the coronary epidemic, an urgent need has emerged for a feasible scientific research [*forscher, lungenarzt* (DE), *scientifique* (FR), *investigadores* (ES), *ricercatori, virologist* (IT)] and Internet browsers, including Facebook, have shown interest in news of scientific discoveries towards a more understanding of the emerging disease and to follow the progress of research and tests to produce the vaccine. Accordingly, the analytical images showed that a wide range of publications touched on this issue through many topics. For example, *chloroquine* was mentioned especially in publications in the French language, due to the controversy caused by the debate about the effectiveness of this drug or not.

As for the economic impact of the Covid-19 epidemic [*euro, industriale* (IT), *euros, compañía* (ES), *firma, شركة* (AR) (DE)], the economic activities of many people around the world have been disrupted as an inevitable consequence of quarantine. This is what led to the emergence of many topics about the axis of work *労働省* (JA), *lavoro* (IT), *workers* (EN), *mitarbeiter* (DE), *travail* (FR), *trabajadores* (ES)], which highlights the busyness of surfers with their livelihoods and their fear of entering into a financial crisis as a result of the disruption of their activities or for fear of losing work in the form of prolonging the crisis for a long time. Accordingly, the term “remote work” [*テレワーク* (JA)] appeared.

And in relation to the users’ interest in the prevention axis [*防止* (JA), *لوقاية* (AR), *protezione* (IT)], the focus was on three measures, namely continuous cleaning *手洗い* (JA), *تغيب* (AR), *hands* (EN), *manos* (ES)] and the use of protective masks *マスク* (JA), *المامات* (AR), *masks* (EN), *maskenpflicht, schutzmasken, masken* (DE), *masque* (FR), *mascherine* (IT)] with social distancing (including self-isolation) [*distancing* (EN), *distanzierung* (DE), *distance* (FR)]. Focusing on these three points in all publications in different languages is an obligation of the user to follow the advice, but it is also a request from the

other to commit to preserving the health of the group that the individual is part of and cannot be active without the participation of others.

As for the axis of communicating with the outside world during the quarantine and the desire to follow all that matters the epidemic as a benefit, many topics expressed the interest of Facebook users in that. Accordingly, we note the existence of a dictionary of news [*ニュース* (JA), *news* (EN)(IT)(FR), *noticias* (ES), *أخبار, أنباء*] and its sources like social networks *ツイート* (JA), *sociales* (ES), *réseaux* (FR)] and others *اذاعة* (AR), *mediaset* (IT)], as users tend to share everything new about the disease. Also, the use of the Internet [*インターネット, データ, オンライン* (JA)] appeared clearly in order to follow the news, entertainment and work from afar.

There is another dimension common to all languages, which is the religious and faith dimensions *الجامع, رمضان, الله* (AR), *jesus, prayer, church* (EN), *papa* (IT)], which we have identified in many subjects. This is mainly due to the feeling of fear and distrust of the future with the large numbers of deaths in many countries such as the United States, Italy, Spain and France. Solidarity [*sociale* (IT), *social, ayudas* (ES), *obdachlose, gibt, hilfe, zusammen* (DE), *الاجتماعي* (AR)] between people also had a share of the extracted subjects, and the intention is to help the needy, from the vulnerable groups, health and financially.

The so-called *rest of the globe* in front of industrial activities largely disrupted, which led to a low level of pollution and improved air quality<sup>26</sup> in countries that followed a general health quarantine for the various productive sectors. This issue appeared on the Italian *inquinamento* language data for April.

### 7.3 Analysis of third period: End April-15 May 2020

By analyzing the data collected until 15 May, as it is illustrated in Fig. 8, we can see the emergence of other issues related to the policy of gradual quarantine lifting (or named Targeted quarantine in Arabic countries such Tunisia) [*déconfinement* (FR), *الحد* (AR), *lockerungen, öffnen, freiheit* (DE), *salir* (ES), *ripartire* (IT)] in many countries that have been affected by the COVID-19 pandemic. This manifests itself by mentioning activities related to the gradual exit of people to their normal lives [*tierpark, maßnahmen, schulen, Kinos* (DE), *colegio, laboral* (ES), *المكّن* (AR), *école, plages, travail* (FR), *market* (EN)], while maintaining the health measures that have

<sup>26</sup><https://www.latribune.fr/opinions/tribunes/quand-les-effets-du-coronavirus-se-voient-depuis-l-espace-841719.html>





**Table 3** Summarizing the discussed topics for each period

|                       |   |
|-----------------------|---|
| January-February 2020 | <ul style="list-style-type: none"> <li>• Acquaintance stage               <ul style="list-style-type: none"> <li>– <i>Virus origin and naming</i>: كورونا (AR), asie, épidemie (FR), cina (IT), Coronavirus (EN) and コロナ (JA), sars and ncov</li> <li>– <i>Emergency affection</i>: allarmismi (IT), urgence (FR), murdered (DE)</li> <li>– <i>Virus outbreak and mobilization of official authorities</i>: الصحة، الغي، السياحة (AR), medico (IT), hausarzt (DE), viajes, vuelve, puerto (ES), クルズ, キャンセル, マスク, キャンセル (JA), masks, hands, quarantine (EN)</li> </ul> </li> </ul>   |
| March-April 2020      | <ul style="list-style-type: none"> <li>• Shock stage and radical change in daily life               <ul style="list-style-type: none"> <li>– <i>A state of alert with the official authorities</i>: 1) スク (JA), emergenza (IT), crisis (ES), crise (FR), ausnahmezustand, coronakrise, gefährlicher (DE), emergenza (IT), fallecido (ES), وفاة (AR), morts, décès, deuil (FR), deaths (EN)</li> <li>– <i>Confrontation of the pandemic by the mobilization of the political and financial capabilities</i>: مكافحة, مجابهة (AR), に対し (JA), lucha, ejército (ES), lutte (FR), 大統領 (JA), sindaco (IT), presidente, gobierno (ES), ministre, gouvernement, autorités (FR), behörde, bund, bundesregierung, regierung (DE), وزير, الحكومة، وزارة (AR), president, government, politics (EN)</li> <li>– <i>Quarantine as taken measure to combat the virus</i>: confinement (FR), الحجر, العزل (AR), quarantäne (DE), cuarente, confinamiento (ES), quarantena (IT), lockdown, quarantine (EN), casa (IT), casa (ES), دارك البيت (AR), house (EN)</li> <li>– <i>Quarantine's related topics</i>:                   <ul style="list-style-type: none"> <li>* <i>Education topic</i>: 子ども, 小学校 (JA), niños (ES), العيال, التربية, جامعة (AR), estudio (ES), écoles, étude (FR), school (EN)</li> </ul> </li> </ul> </li> </ul> |
| March-April 2020      | <ul style="list-style-type: none"> <li>• Shock stage and radical change in daily life               <ul style="list-style-type: none"> <li>* <i>Economic impact of the Covid-19 epidemic</i>: euro, industriale (IT), euros, compãnia (ES), firma (DE), شركة (AR)</li> <li>* <i>Topics about the axis of work</i>: 労働省, テレワーク (JA), lavoro, workers (EN), mitarbeiter (DE), travail (FR), trabajadores (ES)</li> <li>– <i>Prevention Topic</i>: 防止 (JA), للوقاية (AR), 手洗い (JA), تعقيم (AR), hands (EN), manos (ES), マスク (JA), الكمامات (AR), masks (EN), maskenpflicht, schutzmasken, masken (DE), masque (FR), protezione, mascherine (IT), distancing (EN), distanzierung (DE), distance (FR)</li> <li>– <i>Stock exchange of cases and appropriate drug searching</i>: casi (IT), caso (ES), حالة (AR), tests (FR), positiv (DE), positive (EN), positivo (ES), positive (IT), médicaments, traitement (FR), drug (EN), farmaco, terapia (IT), دواء (AR), ワクチン (JA), vaccino (IT), invacuna (ES), impfung, forschler, lungenarzt (DE), vaccine (EN), scientifique (FR), investigadores (ES), ricercatori, virologist (IT)</li> </ul> </li> </ul>  |
| March-April 2020      | <ul style="list-style-type: none"> <li>• Shock stage and radical change in daily life – <i>News and its sources related topic</i>: ニュース, ツイート, ソンターネット, データ, オン (JA), news (EN)(IT)(FR), noticias, sociales (ES), انباء, اخبار (AR), réseaux (FR)               <ul style="list-style-type: none"> <li>– <i>Religious and faith dimensions</i>: رمضان, الله, الجامع (AR), jesus, prayer, church (EN), papa (IT)</li> <li>– <i>Solidarity Topic</i>: sociale (IT), social, ayudas (ES), obdachlose, gibt, hilfe, zusammen (DE), الاجتماعي (AR)</li> </ul> </li> </ul>  |

**Table 3** (continued)

|                        |   |
|------------------------|---|
| End April- 15 May 2020 | <ul style="list-style-type: none"> <li>• Returning to normal life and taking into consideration the coexistence with the virus</li> <li>– <i>Gradual quarantine lifting topic</i>: déconfinement (FR), الحجر (AR), lockerungen, öffnen, freiheit (DE), salir (ES), ripartire (IT), tierpark, maßnahmen, schulen, Kinos (DE), colegio, laboral (ES), المهن (AR), école, plages, travail (FR), market (EN)</li> <li>– <i>Police and law enforcement topic</i>: force (EN), polizia (IT), policia (ES), bußgeldkatalog (DE)</li> </ul> |
|------------------------|---|

## 7.4 Summarizing extracted topics for the whole study period

Table 3 summarizes the similar topics in relation to each period independently from the language that represents one or several countries.

Figure 9 shows the main and common topics discussed in Facebook in relation with different languages. Each topic has been discussed in different language with different probabilities. For instance, *coronavirus* and *covid* are involved as main subject in different languages from the beginning of its appearance to the middle of May 2020 in a cumulative way. As English countries are among the first countries (South Korea) where the infection started early, an important part of English posts, are discussing the quarantine as means of protection. Furthermore, most of French posts are addressing two main topics: health and death topics and this can be explained by the large numbers of deaths that France has reached during this pandemic.

## 7.5 Discussion

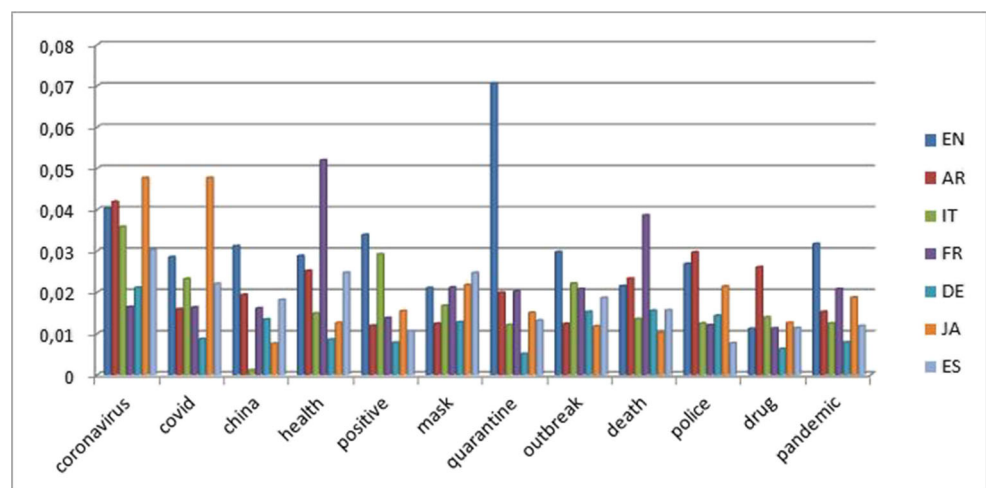
This section is intended to compare what we have reached in this research with previous research in other issues, despite its limited number. [19] work showed how topics can be explored on a more extensive global scale to allow the

public health community understand variations in topics across countries in relation with Zika disease. These data contribute importantly to understanding of disease transmission, disease interventions, and public health communication. This work supports the previously reached social media platforms in following up on pandemics. As in [19], we have succeeded in simulating the reality of living with the COVID-19 pandemic in many languages, representing many countries. However, we did not rely on translating the publications, but rather we used the publications in the language in which they were written. We have also succeeded in creating a new method for presentation that facilitates the process of analyzing results.

Zhang et al. [49] investigated the contents of posts about the Zika virus on Yahoo! Answers, identify and reveal subject patterns about the Zika virus, and analyze the temporal changes of the revealed subject topics over 4 defined periods of the Zika virus outbreak. Multidimensional scaling analysis, temporal analysis, and inferential statistical analysis approaches were used in the study.

The present study shows a result closer to the one about Zika directed by [50] and based on topical analysis concerning what people are tweeting about four disease characteristics and concluding to four main topics (symptoms, transmission, prevention, and treatment) which are existing also for COVID-19.

**Fig. 9** Analysis of common topics in relation to their normalized weights in different languages



## 8 Conclusion and Future Works

This concurrent research with the mobilization of research capabilities across all countries on the COVID-19 pandemic is considered the first of its kind, which depends on the analysis of the data of the social platform Facebook, which is the first platform globally through the number of active users. This research work focused on multilingualism in seven languages English (EN), Arabic (AR), Italian (IT), Spanish (ES), French (FR), German (DE) and Japanese (JA), representing a large group of countries in order to analyze the cognitive development of people about COVID-19 and the development of topics raised over time from 1 January 2020 to 15 May 2020.

The results of analyzing the data collected and analyzed by adopting three cumulative periods as the first extends from January to February, and the second from January to April with the allocation of the months of March and April each with a graph, while the third period extends 15 May.

This work highlighted the ability to simulate the reality of life through the time of COVID-19 in various countries of the world and in many languages. It can have an outlook for emerging issues. This proposed approach also remains valid for exploring other events and gives the possibility for an in-depth analysis of well-selected topics in recursive way. In fact, the LDA output is considered as fuzzy classification affecting the posts to the extracted topics. Moreover, the presentation module is the first attempt for expressing topics in such graph structure providing a meaningful way to interpret and spot main issues such COVID-19 related research.

The topics revolved around several axes, which began to feel the seriousness of the coming corona crisis, with a weakness in the number of publications regarding the COVID-19 virus during January and February, except for the English language as a result of its global use and Japanese due to the early spread of the disease in Japan and East Asian countries in general. The beginnings were characterized by focusing on the origin of the virus with the circulation of several terms related to the virus and everything related to the health sector in the concerned countries. In the second phase, which reaches the months of March and April, the Facebook users of along this target research discussed, in various degrees, everything related to mobilizing countries to confront the pandemic, while mentioning the procedures followed for protection that remain among the most important of the quarantine and the effects that accompanied it on many levels, including the family, social, economic and the politician. It also touched on the death issue as a result of the high number of victims of this disease. People also showed interest in scientific research through following up on everything new about treatment, especially the discovery of vaccine. In a third stage that extends to the middle of May, we have noticed a

gradual emergence of people's desire to lift the quarantine, even if only partially, and their desire to gradually return to normal life appeared with everyone's desire to follow health procedures from everyone in order to preserve the safety of the participants in the public space.

A summary is given in order to show the evolution of the topics according to the specific language, the period and a normalized weight provided by LDA.

This work, as it depicted the reality experienced by the world in a unique experience with the COVID-19 pandemic, can be proactive in order to sense the changes that are in the process of occurring or that may occur. It is also worth noting that this system can be reproduced to work on any previous or later issue in order to have a deeper study in relation with the temporal dimension through the publications of social media platforms, especially Facebook.

In future works, we plan to broaden our work to cover other languages and to go deeper in the analysis by developing a recursive process able to zoom in the topics by extracting the sub-topics and building predictive models which is favored by the probabilistic generative of LDA. Moreover, we will focus on the integration of semantic technologies to analyse the semantic dimension with embedding models in the measurement of the multilingual and cross-lingual semantic similarity/relatedness.

**Acknowledgment** The work was supported by the Ministry of Higher Education and Scientific Research of Tunisia (MHESR) in the framework of Federated Research Project PRFCOV19-D1-P1.

## References

1. Sebei H, Taieb MAH, Aouicha MB (2018) Review of social media analytics process and big data pipeline, *Social Netw Anal Mining* 8:30:1–30:28
2. Teodorescu H-N (2015) Using analytics and social media for monitoring and mitigation of social disasters, *Procedia Engineering* 107:325–334
3. Joseph JK, Dev KA, Pradeepkumar A, Mohan M (2018) Big data analytics and social media in disaster management. In: *Integrating Disaster Science and Management*, Elsevier, pp 287–294
4. Landwehr PM, Carley KM (2014) *Social Media in Disaster Relief*, Springer Berlin Heidelberg, Berlin, Heidelberg, 225–257
5. Doan S, Vo B-KH, Collier N (2012) An analysis of twitter messages in the 2011 tohoku earthquake. In: P Kostkova, M Szomszor, D Fowler (Eds.), *Electronic Healthcare*, Springer Berlin Heidelberg, Berlin, Heidelberg, pp 58–66
6. Miyabe M, Miura A, Aramaki E (2012) Use trend analysis of twitter after the great east japan earthquake. In: *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work Companion, CSCW '12*, Association for Computing Machinery, New York, NY, USA, pp 175–178
7. Sakaki T, Okazaki M, Matsuo Y (2010) Earthquake shakes twitter users: Real-time event detection by social sensors. In: *Proceedings of the 19th International Conference on World Wide Web, WWW '10*, Association for Computing Machinery, New York, NY, USA, pp 851–860

8. PEARY B, Shaw R, TAKEUCHI Y (2012) Utilization of social media in the east japan earthquake and tsunami and its effectiveness. *Journal of Natural Disaster Science* 34:3–18
9. Daga RRM (2017) Social network analysis of tweets on typhoon during haiyan and hagupit. In: *Proceedings of the 8th International Conference on Computer Modeling and Simulation, ICCMS '17*, Association for Computing Machinery, New York, NY, USA, pp 151–154
10. Ulvi O, Lippincott N, Khan MH, Mehal P, Bass M, Lambert K, Lentz E, Haque U (2019) The role of social and mainstream media during storms. *Journal of public health and emergency*, 3
11. Kankanamge N, Yigitcanlar T, Goonetilleke A, Kamruzzaman M (2020) Determining disaster severity through social media analysis: Testing the methodology with south east queensland flood tweets. *International journal of disaster risk reduction*, 42
12. Ahmed W (2018) Using twitter data to provide qualitative insights into pandemics and epidemics
13. Fan B, Fan W, Smith C, Garner HS (2020) Adverse drug event detection and extraction from open data: A deep learning approach. *Information Processing and Management* 57:102–131
14. Pizzuti AG, Patel KH, McCreary EK, Heil E, Bland CM, Chinaeke E, Love BL, Bookstaver PB (2020) Healthcare practitioners' views of social media as an educational resource. *PLOS ONE* 15:1–16
15. Ding H, Zhang J (2010) Social media and participatory risk communication during the h1n1 flu epidemic: A comparative study. *China Media Research* 6:80–91
16. Achrekar H, Gandhe A, Lazarus R, Yu S-H, Liu B (2011) Predicting flu trends using twitter data. In: *2011 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPs)*, pp 702–707
17. Lee K, Agrawal A, Choudhary A (2017) Forecasting influenza levels using real-time social media streams. In: *2017 IEEE International Conference on Healthcare Informatics (ICHI)*, pp 409–414
18. Sharma M, Yadav K, Yadav N, Ferdinand KC (2017) Zika virus pandemic—analysis of facebook as a social media health information platform. *American Journal of Infection Control* 45:301–302
19. Pruss D, Fujinuma Y, Daughton A, Paul M, Arnot B, Szafir D, Boyd-Graber J (2019) Zika discourse in the americas: A multilingual topic analysis of Twitter. *PlosOne*
20. Zarrad A, Jaloud A, Alsmadi I (2014) The evaluation of the public opinion - a case study: Mers-cov infection virus in ksa. In: *2014 IEEE/ACM 7th International Conference on Utility and Cloud Computing*, pp 664–670
21. Tran T, Lee K (2016) Understanding citizen reactions and ebola-related information propagation on social media. In: *2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pp 106–111
22. Missier P, Romanovsky A, Miu T, Pal A, Daniilakis M, Garcia A, Cedrim D, da Silva Sousa L (2016) Tracking dengue epidemics using twitter content classification and topic modelling. In: *Current Trends in Web Engineering - ICWE 2016 International Workshops, DUI, TELERISE, SoWeMine, and Liquid Web, Lugano, Switzerland, June 6-9, 2016, Revised Selected Papers*, pp 80–92
23. Sicilia R, Giudice SL, Pei Y, Pechenizkiy M, Soda P (2018) Twitter rumour detection in the health domain. *Expert Systems with Applications* 110:33–40
24. Alshaabi T, Arnold MV, Minot JR, Adams JL, Dewhurst DR, Reagan AJ, Muhamad R, Danforth CM, Dodds PS (2020) How the world's collective attention is being paid to a pandemic: COVID-19 related 1-gram time series for 24 languages on Twitter
25. Barkur G, Prabhu V, Kamath G (2020) Sentiment analysis of nationwide lockdown due to covid 19 outbreak: Evidence from india. *Asian Journal of Psychiatry* 51:102–089
26. Chen Q, Min C, Zhang W, Wang G, Ma X, Evans R (2020) Unpacking the black box: how to promote citizen engagement through government social media during the covid-19 crisis. *Computers in human behavior*. <https://doi.org/10.1016/j.chb.2020.106380>
27. Li L, Zhang Q, Wang X, Zhang J, Wang T, Gao T-L, Duan W, Tsoi K, Wang F-Y (2020) Characterizing the propagation of situational information in social media during covid-19 epidemic: A case study on weibo. *IEEE Transactions on Computational Social Systems* PP:1–7
28. Limaye R, Sauer M, Ali J, Bernstein J, Wahl B, Barnhill A, Labrique A (2020) Building trust while influencing online covid-19 content in the social media world, the lancet digital health
29. Zhou C, Su F, Pei T, Zhang A, Du Y, Luo B, Cao Z, Wang J, Yuan W, Zhu Y, Song C, Chen J, Xu J, Li F, Ma T, Jiang L, Yan F, Yi J, Hu Y, Xiao H (2020) Covid-19: Challenges to gis with big data, geography and sustainability
30. Chen E, Lerman K, Ferrara E (2020) Covid-19: The first public coronavirus twitter dataset. *arXiv:2003.07372*
31. Alqurashi S, Alhindi A, Alanazi E (2020) Large arabic twitter dataset on covid-19. *arXiv:2004.04315*
32. Banda JM, Tekumalla R, Wang G, Yu J, Liu T, Ding Y, Chowell G (2020) A large-scale covid-19 twitter chatter dataset for open scientific research—an international collaboration. *arXiv:2004.03688*
33. Boberg S, Quandt T, Schatto-Eckrodt T, Frischlich L (2020) Pandemic populism: Facebook pages of alternative news media and the corona crisis - a computational content analysis
34. Cinelli M, Quattrociocchi W, Galeazzi A, Valensise CM, Brugnoli E, Schmidt AL, Zola P, Zollo F, Scala A (2020) The COVID-19 social media infodemic. *arXiv:2003.05004*
35. Gao Z, Yada S, Wakamiya S, Aramaki E (2020) NAIST COVID: multilingual COVID-19 twitter and weibo dataset. *arXiv:2004.08145*
36. Kleinberg B, van der Vegt I, Mozes M (2020) Measuring emotions in the COVID-19 real world worry dataset. *arXiv:2004.04225*
37. Kuchler T, Russel D, Stroebel J (2020) The geographic spread of covid-19 correlates with structure of social networks as measured by facebook, technical report, national bureau of economic research
38. Lopez CE, Vasu M, Gallemore C (2020) Understanding the perception of COVID-19 policies by mining a multilanguage twitter dataset. *arXiv:2003.10359*
39. Zarei K, Farahbakhsh R, Crespi N, Tyson G (2020) A first instagram dataset on covid-19
40. Perrotta D, Grow A, Rampazzo F, Cimentada J, Del Fava E, Gil-Clavel S, Zagheni E (2020) Behaviors and attitudes in response to the covid-19 pandemic: Insights from a cross-national facebook survey. *medRxiv*
41. Jelodar H, Wang Y, Yuan C, Feng X, Jiang X, Li Y, Zhao L (2019) Latent dirichlet allocation (lda) and topic modeling: Models, applications, a survey. *Multimedia Tools Appl* 78:15169–15211
42. Blei DM, Ng AY, Jordan MI (2003) Latent dirichlet allocation. *J Mach Learn Res* 3:993–1022
43. Follett L, Geletta S, Laugerman M (2019) Quantifying risk associated with clinical trial termination: A text mining approach. *Inf Process Manag* 56:516–525
44. Liu L, Tang L, Dong W, Yao S, Zhou W (2016) An overview of topic modeling and its current applications in bioinformatics. *SpringerPlus* 5:1608

45. Damevski K, Chen H, Shepherd DC, Kraft NA, Pollock LL (2018) Predicting future developer behavior in the IDE using topic models. *IEEE Trans Software Eng* 44:1100–1111
46. Amara A, Taieb MAH, Aouicha MB (2017) Identifying i-bridge across online social networks. In: 14th IEEE/ACS International Conference on Computer Systems and Applications, AICCSA 2017, Hammamet, Tunisia, October 30 - Nov. 3, 2017, pp 515–520
47. McCallum AK (2002) Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu>
48. Missier P, McClean C, Carlton J, Cedrim D, Silva L, Garcia A, Plastino A, Romanovsky A (2017) Recruiting from the network: Discovering twitter users who can help combat zika epidemics. In: J Cabot, R De Virgilio, R Torlone (Eds.), *Web Engineering*, Springer International Publishing, Cham, pp 437–445
49. Zhang J, Chen Y, Zhao Y, Wolfram D, Ma F (2020) Public health and social media: A study of Zika virus-related posts on Yahoo! Answers. *Journal of the Association for Information Science & Technology* 71:282–299
50. Miller M, Banerjee T, Muppalla R, Romine W, Sheth PA (2017) What are people tweeting about zika? an exploratory study concerning its symptoms, treatment, transmission, and prevention, *jmir public health and surveillance*

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.