

---

**This report is not subject to GEROS assessment.**

It is an explorative evaluative report and is not written in line with the evaluation report standards.

# **Real-Time Assessment of UNICEF's Ongoing Response to COVID-19 in Europe and Central Asia**

Social Media Listening (SML) Exploration

Denis Nikitin, Umer Naeem, and Paul Jasper

April 2022

## About Oxford Policy Management

Oxford Policy Management (OPM) is committed to helping low- and middle-income countries achieve growth and reduce poverty and disadvantage through public policy reform.

We seek to bring about lasting positive change using analytical and practical policy expertise. Through our global network of offices, we work in partnership with national decision makers to research, design, implement, and evaluate impactful public policy.

We work in all areas of social and economic policy and governance, including health, finance, education, climate change, and public sector management. We draw on our local and international sector experts to provide the very best evidence-based support.

The contact point for the client is Paul Jasper ([paul.jasper@opml.co.uk](mailto:paul.jasper@opml.co.uk)). The contact point at UNICEF is Saltanat Rasulova ([srasulova@unicef.org](mailto:srasulova@unicef.org)).

## Table of contents

List of tables and figures .....	iii
List of abbreviations .....	v
1 Introduction .....	1
1.1 Background.....	1
1.2 Objectives.....	2
1.3 Summary of work phases.....	2
2 Methods .....	4
2.1 Obtaining and using data from UNICEF’s TW platform .....	4
2.2 The main queries on which this analysis is based .....	6
2.3 Text pre-processing prior to analysis.....	7
2.4 Data analysis .....	7
3 Results .....	9
3.1 Education.....	9
3.2 SP.....	19
4 Conclusions, lessons learned, and recommendations .....	34
4.1 Conclusions .....	34
4.2 Lessons learned.....	37
4.3 Methodological recommendations for further use of SML.....	38
Annex A: SML memo .....	40
Annex B: Search queries .....	46
Annex C: Topic modelling for Albania and Tajikistan.....	49

## List of tables and figures

Table 1: Number of results in each query.....	7
Table 2: Query 1 results by country.....	9
Table 3: Education—topic modelling.....	13
Table 4: Mentions of UNICEF in SP discourse.....	25
Table 5: Albania—word clouds of top 70 terms central to all SP discourse and SP discourse mentioning UNICEF by period (from left to right, Periods 1, 2, and 3).....	26
Table 6: Albania—topic modelling results.....	27
Table 7: Tajikistan—topic modelling results.....	32
Figure 1: Outline of work phases.....	3
Figure 2: Example of a ‘snippet’ download from TW.....	4
Figure 3: Example of data filtering to ‘complete’ text.....	5
Figure 4: Text analysis workflow.....	6
Figure 5: Social media source in the first education query.....	10
Figure 6 Social media posts in the first education query over time.....	11
Figure 7: Education—absolute wordcount for words with five or more mentions.....	12
Figure 8: Education—relative word frequency by country group.....	12
Figure 9: Education—number of words that explain 60% of each topic.....	13
Figure 10: Education—sentiment of topics.....	14
Figure 11: Education—topics by country group.....	15
Figure 12: Education—topic importance over time.....	15
Figure 13: Education—word cloud relating to Topic 3.....	16
Figure 14: Education—word cloud relating to Topic 5.....	17
Figure 15: Proportion of UNICEF mentions in non-UNICEF query on TW.....	18
Figure 16: Education—sentiment in UNICEF-related posts on the TW platform.....	18
Figure 17: Total number of occurrences of SP posts in select countries (pre-filtering).....	21
Figure 18: Frequency of mentions of SP (population-adjusted) and relative size of SP sector (% of GDP).....	22
Figure 19: Evolutions in SP mentions in Albania by period.....	22
Figure 20: Evolution of SP mentions in Tajikistan by period.....	23
Figure 21: The 20 most frequent words in Albania in each period (normalised frequency) ..	24

Figure 22: Albania—change in centrality of select word clusters (word frequencies normalised by occurrence of the most common word in the period) ..... 25

Figure 23: Sentiment and emotion profile of the Albania’s SP discourse in three study periods (scoring based on the NRC lexicon) ..... 29

Figure 24: The 20 most frequent words in Tajikistan in each period (normalised frequency)30

Figure 25: Tajikistan—change in centrality of select words (word frequencies normalised by occurrence of the most common word in the period) ..... 30

Figure 26: Tajikistan—word clouds of up to top 70 terms central to all SP discourse and SP discourse mentioning UNICEF ..... 31

Figure 27: Sentiment and emotion profile of Tajikistan’s SP discourse in three study periods (scoring based on the NRC lexicon) ..... 33

Figure 28: OPM’s proposed process workflow to conduct data analysis ..... 43

Figure 29: Normalised frequency distributions of top 20 words in the UNICEF-related discourse in Albania ..... 49

Figure 30: SP topic modelling Albania—Period 1 ..... 51

Figure 31: SP topic modelling Albania—Period 2 ..... 52

Figure 32: SP topic modelling Albania—Period 3 ..... 53

Figure 33: SP topic modelling Tajikistan—Period 1 ..... 54

Figure 34: SP topic modelling Tajikistan—Period 2 ..... 55

Figure 35: SP topic modelling Tajikistan—Period 3 ..... 56

## List of abbreviations

ECAR	Eastern Europe and Central Asia Region
ECARO	Europe and Central Asia Region Office
IFRC	International Federation of the Red Cross and Red Crescent
LDA	Latent Dirichlet Allocation
NLP	Natural Language Processing
NRC	National Research Council Canada
OPM	Oxford Policy Management
RTA	Real-Time Assessment
SML	Social Media Listening
SP	Social Protection
TW	TalkWalker
UNFPA	United Nations Population Fund
UNICEF	United Nations Children's Fund

# 1 Introduction

## 1.1 Background

This report forms part of the larger Real-Time Assessment (RTA) of the response of the United Nations Children's Fund (UNICEF) to COVID-19 in the Eastern Europe and Central Asia Region (ECAR). To date, the RTA exercise has consisted of two rounds: RTA-1 and RTA-2. RTA-1 focused on the earlier stage of the COVID-19 response between March and October/November 2020 and aimed to give a broad view of UNICEF's response to the pandemic in ECAR across all areas of UNICEF country office's response.

RTA-2 aims to provide an in-depth understanding of UNICEF's pandemic response in two areas: **education and social protection (SP)**. It focuses on nine Europe and Central Asia Regional Office (ECARO) countries that indicated their interest in participating in this study to the office: Albania, Azerbaijan, Bosnia & Herzegovina, Montenegro, North Macedonia, Serbia, Tajikistan, Turkey, and Uzbekistan. The choice of these two deep-dive areas—education and SP—was informed by the preferences of individual country offices, expressed through a survey, and the findings of the Round 1 report. This choice was subsequently endorsed by the ECAR Deputy Regional Director.

Part of the work of both RTA-1 and RTA-2 was a so-called natural language processing (NLP) component. In RTA-1, the NLP work focused on analysing a static text corpus of documents provided by ECARO that detailed the UNICEF response to the ongoing pandemic in ECAR. In this phase (RTA-2), the focus was to switch from an analysis of this static corpus to a more 'real-time' analysis of text derived from the internet and—more specifically—social media. This explains the subtitle to this report: a social media listening (SML) exploration.

During the inception phase to RTA-2, the research team, together with UNICEF, decided to make use of the existing access to TalkWalker (TW) for the purposes of this SML exploration. TW<sup>1</sup> is an online application that provides users with the possibility of querying a wide range of online media and social media sources and analysing the resulting data on the platform itself. All the SML analyses implemented in the context of RTA-2 and presented in this report were done with data queried on the TW platform. As explained in Section 1.3, the researchers started this work by exploring the usage and functionalities of the TW platform.

The remainder of this report is structured as follows. In this section, we delineate the objectives of this workstream and describe the process implemented to achieve these objectives. In Section 2, we summarise the methods employed to both obtain and analyse the social media data. In Section 3, we present results. Section 3.1 focuses on our analysis relating to the topic area of education, while Section **Error! Reference source not found.** focuses on SP. We conclude in Section 4 with a description of conclusions, lessons learned, and recommendations for the way forward.

---

<sup>1</sup> <https://www.talkwalker.com/>.

## 1.2 Objectives

As per the inception report for the RTA-2, the SML workstream planned to investigate the following questions.

- What are the key issues related to SP/education that are being discussed on social media? Is the sentiment behind these discussions positive, negative, or neutral?
- How have the topics in SP and education discussed online, and the sentiments behind them, changed over time?
- To what extent does the programming and support provided by UNICEF align with the key topics discussed in the social media? Are mentions of UNICEF programmes or support in SP and educations associated with positive or negative sentiments?

In terms of outputs, two objectives of the SML workstream were to deliver:

- a) a methodological note describing the work that has been carried out using the RTA dashboard deployed on the TW platform to extract insights related to RTA questions, any challenges faced, and lessons learned; and
- b) a sample analytical report summarising the findings of the NLP analysis as it pertains to up to two RTA questions for one sample country in each of the deep-dive areas.

This report represents the report mentioned under b). Note, however, that we expanded the list of countries covered in this work, as explained in Section 2. Responding to point a), a methodological memo was presented to UNICEF in October 2021 and is included in this report in Annex A. It is important to note, however, that an analysis was not carried out on TW, but separately using statistical programming (see Section 2.1) for detail. This report updates the methodological description included in that memo. Finally, it is important to mention that we were not able to develop a fully functional automated dashboard within the context of RTA-2 for reasons explained in the methods section of this report and in Annex A. However, all analysis implemented in the context of this SML exploration was implemented using statistical programming software with HTML output that can in future easily be integrated into a dashboard. In fact, at the time of finalising this report, the Oxford Policy Management (OPM) team has shared this code with the UNICEF Evaluation Office for further integration into analytical processes at UNICEF.<sup>2</sup>

## 1.3 Summary of work phases

As described in our inception report, the NLP and SML work for this RTA-2 phase was exploratory in character. A considerable amount of effort, therefore, was spent on interacting with the existing SML platform from UNICEF (TW) and on understanding how data derived from this platform could be used for the purpose of evaluation. Figure 1 presents the phases in which the SML work for this RTA-2 was implemented. These can be described as follows.

- Between August and September 2021, the team mainly explored the usage of the TW platform and its functionality. This included frequent interaction with the TW technical team (via email and telephone calls) and trial runs for queries and analyses on the platform.

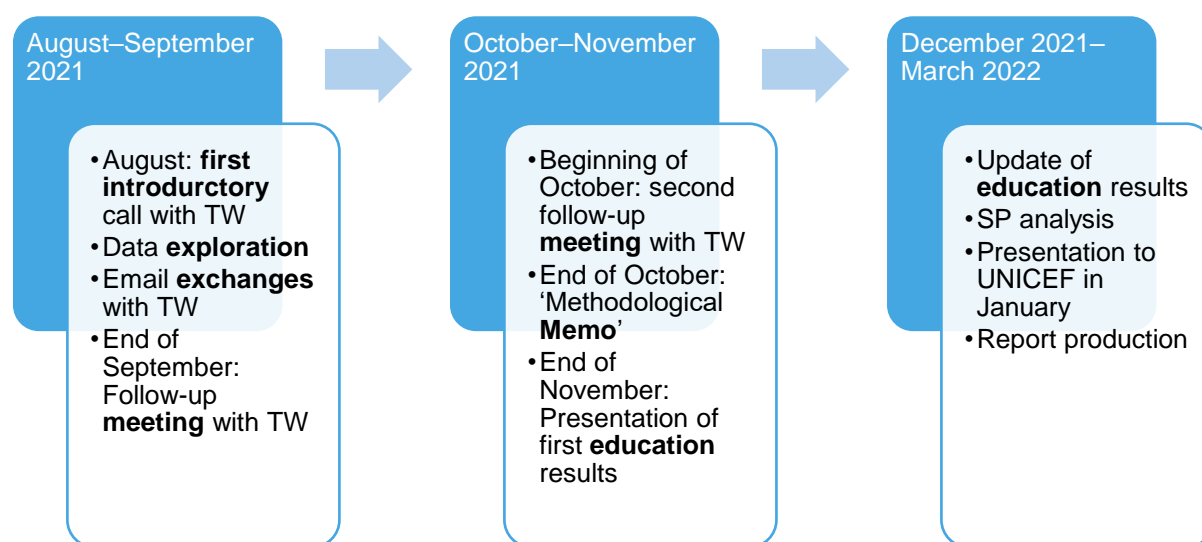
---

<sup>2</sup> The R and R-Markdown codes were shared on 11 February 2022 with the intention of integrating this into dashboard functionality at UNICEF.



- Between October and November 2021, the team developed the proposed methodology for analysis in the context of this RTA-2. This was presented in a memo to UNICEF at the end of October 2021. In November, a first round of analyses on the topic of education were implemented and presented to UNICEF.
- Between December 2021 and January 2022, the team updated this education analysis and implemented the SP analysis. This SP analysis was presented to UNICEF in January. Results were summarised in the present report.

**Figure 1: Outline of work phases**



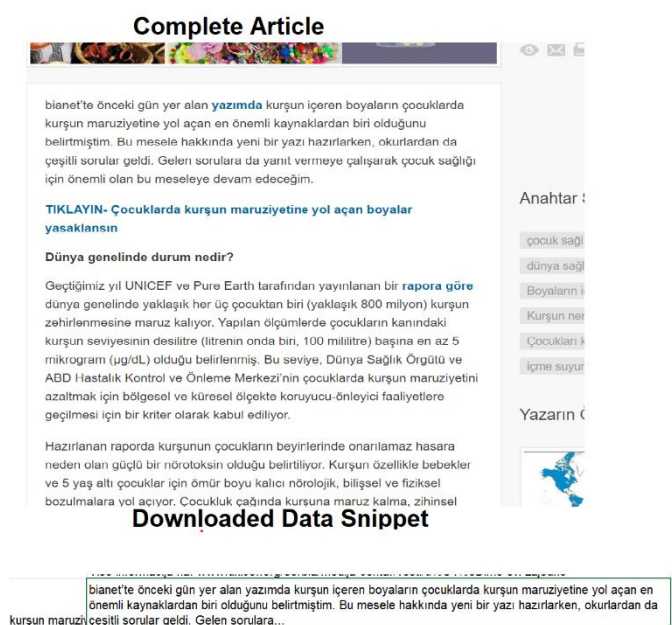
## 2 Methods

### 2.1 Obtaining and using data from UNICEF's TW platform

As mentioned above, in the inception phase to this project, the research team, together with UNICEF, decided to use the TW platform to obtain data to be used in this SML analysis. The original idea was to use this platform to define queries, obtain data, and analyse them using the on-app dashboard functionality. However, in the first work phase described in Figure 1, the research team encountered several issues that prevented us from proceeding with this plan. We present these in more detail in Annex A. They can be summarised as follows.

- **Limited pre-processing:** The pre-processing functionality on the TW platform is limited. While some pre-processing is feasible (e.g. removing hashtags), this only works manually (i.e. by clicking on certain words). This means that implementation would be very laborious and limit replicability in the longer term.
- **Limited uploading functionality:** To address the pre-processing issue, an alternative approach could have been to run the query on TW, download the data, pre-process them, and then reupload them onto the platform for analyses. However, the TW platform did not handle reuploads of processed data well, resulting in errors. Despite several trial runs by the research team, we were thus not able to implement this reupload within the timeframe of this project.

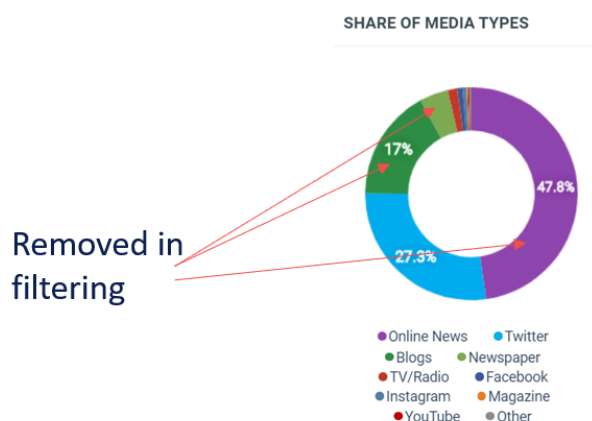
**Figure 2: Example of a 'snippet' download from TW**



- **Incomplete data:** The most significant issue we encountered was that the data derived from TW queries were incomplete. More specifically, this referred to the text relating to each result derived from TW queries. If the text for these results was behind a paywall (i.e. not publicly available), the results used by TW were in fact only a 'snippet' of the full text. More significantly, when downloading data, long texts from

publicly available sources (e.g. a newspaper article) would also be shortened by TW into 'snippets' (see Figure 2 for an example of this from a Turkish article). In fact, according to TW, the full text can only be seen and exported from 11 social media sources.<sup>3</sup> For the current project, this meant limiting our analysis to data from these sources. Figure 3 shows an example of what this means in the context of an education query. In essence, the data used for the analysis are limited to Twitter posts.

**Figure 3: Example of data filtering to 'complete' text**



- Translation:** Most of the text analysed in the context of this project was originally not posted in English but in other languages or even other scripts (e.g. Cyrillic script). Translation into English was therefore necessary. While TW does provide translation on the platform, bulk translation for downloads would have been an addition the project would have had to pay for. All data were therefore downloaded in the original text and then translated in bulk using a freeware Google Sheets function that employs Google Translate.<sup>4</sup> This also meant we had to deal with hashtags and symbols (emojis) that people included in their posts prior to translation.
- Downloading data:** Because the analysis carried out for the purposes of this assessment was historical (i.e. observing posts created throughout the COVID-19 period since the beginning of 2020), the queries that were run provided a large set of results. It turned out downloading data from TW could only be done in day-by-day chunks of 50,000 results, which meant that downloading the full set of some queries took several days.

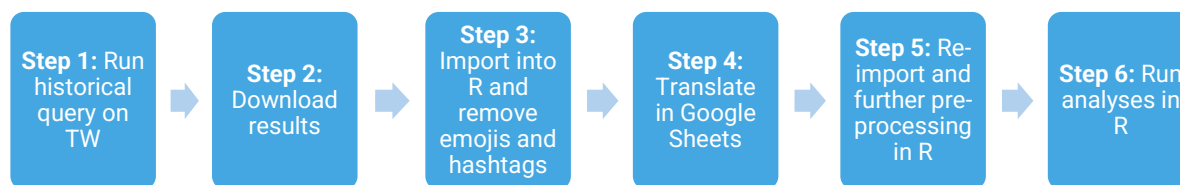
These issues and the ways in which the research team decided to tackle them led to a workflow summarised in Figure 4. As a first step, a query was run on TW to obtain results relevant to the themes and countries looked at in the context of this analysis (education and SP). As a second step, these results were downloaded, which could sometimes take several days. As a third step, an initial pre-processing was implemented in that emojis and hashtags

<sup>3</sup> In an email sent in November 2021, TW confirmed that the sources for which full content can be seen are Facebook, Twitter, YouTube, LinkedIn, Instagram, Mixcloud, Soundcloud, Vimeo, Dailymotion, VKontakte, and Twitch.

<sup>4</sup> See <https://support.google.com/docs/answer/3093331> for more information.

were removed from the data. The resulting text corpus was then translated into English using Google Sheets. A fifth step consisted of further pre-processing (explained in more detail in Section 2.3). Finally, as a sixth step, the data were analysed.

**Figure 4: Text analysis workflow**



## 2.2 The main queries on which this analysis is based

As described in Section 1.2, the SML analysis focused on two key topics: SP and education. Similarly, the analysis was supposed to focus on a subset of the countries in the ECAR region. The queries used on TW to obtain social media data that could then be analysed were derived from these requirements. The full queries can be found in Annex B. Here we summarise their content.

- **Query 1—Education, UNICEF, COVID-19:** A first query focused on obtaining posts that included terms relating to the theme of education, but which also included mentions of UNICEF and COVID. This query covered the four countries included in the education deep-dive report: Azerbaijan, Turkey, Serbia, and Bosnia & Herzegovina. Data were obtained on this query between December 2019 and December 2021.
- **Query 2—Education, COVID-19:** A second query focused on the same countries and themes as in Query 1. However, UNICEF was removed from this search query. The objective was to understand, in general, what the conversation on education and COVID-19 was in the same focus countries as in Query 1 and to see whether UNICEF was being mentioned ‘naturally’. Data here were obtained for the period between January 2020 and January 2022, the two-year period that we could cover prior to the first running of this query. Note that, because of the large scope of this query, the sources were limited to those that TW identified as ‘full text’ sources (Section 2.1). Also note that the data from this query were analysed only briefly, with a view of understanding the proportion of UNICEF mentions.
- **Query 3—SP:** A third query focused on issues related to SP. The analysis of the SP queried data were implemented after the education analysis and followed a slightly different analytical strategy. We describe the query in more detail in Section 3.2.1, but in essence the objective was to obtain data from online discourse relating to issues that dealt with ‘social’ and ‘protection’ as expressed by different words that might reflect these dimensions and as validated by experts in this area from the countries this query focused on: Albania and Tajikistan. The SP query covered the period between 10 January 2020 and 26 December 2021.

Table 1 presents a summary of the number of results obtained by implementing these queries. The pre-filtering row refers to the results obtained prior to the filtering to 'full text' sources. The post-filtering row refers to the results that remained after filtering out the snippets. These are thus the data that the main analyses (for which results are presented in Section 3) are based on.

**Table 1: Number of results in each query**

	Query 1	Query 2	Query 3
Pre-filtering	~25,100	NA	210,775
Post-filtering	7,287	461,839	7,787

## 2.3 Text pre-processing prior to analysis

As described in Section 2.1 and summarised in Figure 4, apart from filtering, the text data obtained were pre-processed prior to analysis. Two 'general' pre-processing steps—removing emojis, hashtags, and translation—resulted in a text corpus in English that could then be read into R for NLP proper, from Step 5 in Figure 4 onwards. The main pre-processing implemented in Step 5 can be summarised as follows:

- first, retweets and mentions were removed;
- second, URLs, numbers, and punctuation were removed;
- third, stop words were removed;<sup>5</sup> and
- fourth, for network analyses and word clouds, stemming was also implemented.<sup>6</sup>

All pre-processing was implemented in R software using the Text Mining Package.<sup>7</sup>

## 2.4 Data analysis

The resulting corpus of data was then analysed in R, mostly using the 'quanteda' and 'topicmodels' packages.<sup>8</sup> The exact types of analysis implemented varied by themes (described further below) but, in general, the analyses focused on the following approaches, implemented one after the other:

- **descriptive analysis of posts:** this meant looking at results (i.e. posts) by country, authors, sources, and over time;

<sup>5</sup> A list of stop words can be found here: <https://rdrr.io/rforge/tm/man/stopwords.html>.

<sup>6</sup> 'Stemming' is a process that reduces words to their 'base components', such as 'win' in the case of 'winning, winner, wins'. See [www.rdocumentation.org/packages/RTextTools/versions/1.4.3/topics/wordStem](http://www.rdocumentation.org/packages/RTextTools/versions/1.4.3/topics/wordStem) for the R documentation and <https://snowballstem.org/> for an explanation of the algorithm.

<sup>7</sup> This was first presented here: [www.jstatsoft.org/article/view/v025i05](http://www.jstatsoft.org/article/view/v025i05).

<sup>8</sup> See <https://quanteda.io/> and <https://cran.r-project.org/web/packages/topicmodels/topicmodels.pdf> for more information.

- **description of words in context:** this meant looking at how and in which patterns words occurred and co-occurred—for example, we counted the number of times certain words appeared, which words appeared most frequently, and how this varied by characteristics of results. It also involved looking at word clouds and networks of words;
- **sentiment analysis** using a dictionary-based approach in which words are assigned sentiment scores based on a dictionary pre-coded by humans—in the present case, we employed the sentiment lexica available via the R package ‘syuzhet’,<sup>9</sup> which allowed us to provide an analysis of positive versus negative sentiments and a set of eight predefined emotions; and
- **topic modelling:** As a final analytical tool, we deployed topic modelling to the text corpora. Topic modelling assumes that documents have been written with certain topics in mind that shape the text that forms documents. These topics are not necessarily explicitly revealed in the text, but reveal themselves when reading documents and corpora of documents. This can be easy for humans, and the objective of topic modelling in NLP is to train algorithms to do this automatically on large sets of text. In the present case, we used the R package ‘topicmodels’ on the results from our queries.<sup>10</sup> More specifically, we implemented a Latent Dirichlet Allocation (LDA) model, which assumes there is a pre-specified, fixed set of topics in a text corpus that can be represented by clusters of words that frequently appear together. In our analyses, we pre-specified the number of topics to look for using topic modelling and reviewed the clusters of words we obtained after modelling. Labelling these topics then has to be carried out by the research team.

It is important to emphasise here that the analytical process varied slightly between the education and SP topic areas. As presented in Figure 1, the education topic area was the first to be explored in the context of this project. This meant that, for the SP analysis, we could build on some lessons learned from this first exploration. In fact, as can be seen in Section **Error! Reference source not found.**, we used the same analytical methods listed above to explore the SP discourse online. However, we used the results slightly differently (e.g. by looking at relative normalised word frequencies) and split the analysis into different periods relating to how the COVID-19 pandemic unfolded in the two countries of interest.

---

<sup>9</sup> See the documentation at [www.rdocumentation.org/packages/syuzhet/versions/1.0.6](http://www.rdocumentation.org/packages/syuzhet/versions/1.0.6) for more information.

<sup>10</sup> See the documentation at <https://cran.r-project.org/web/packages/topicmodels/> for more information.

## 3 Results

In this section, we present a selection of key results of the analyses implemented on the text corpora that derive from the queries and the processing steps described in Section 2.2 and 2.3. All analyses were implemented in the statistical programming language R. The complete set of results is presented in HTML files that were produced in R Markdown, and which will be shared with UNICEF separately.

### 3.1 Education

As discussed in Section 2.2, data resulting from two queries were analysed to provide insights into social media posts relating to education in Azerbaijan, Turkey, Serbia, and Bosnia & Herzegovina. We present an analysis of Query 1 first in Section 3.1.1, and a summary of key findings on Query 2 in Section 3.1.2.

#### 3.1.1 Query 1—Education, UNICEF, and COVID-19

##### Descriptive analysis of posts

As shown in Table 1, this query, which covered the four education deep-dive countries for the period between December 2019 and December 2021, yielded a total of 7,287 results after downloading and pre-processing the data. Table 2 shows how these results are distributed across the countries that we looked at. By far the largest proportion of results came from Turkey, followed by Azerbaijan.

**Table 2: Query 1 results by country**

Country	Number of results	Proportion (%)
Azerbaijan	1,061	14.6%
Bosnia & Herzegovina	140	1.9%
Serbia	198	2.7%
Turkey	5,888	80.8%
<b>Total</b>	<b>7,287</b>	<b>100.0%</b>

Figure 5 shows how the sources from which these results derive vary by country. For this graph, Serbia and Bosnia & Herzegovina were grouped together in the country group 'Balkans' because of their low absolute volume in the text corpus at hand (see Table 2). Overall, about 95% of all results we are looking at are posts on Twitter. This high proportion holds across the different country locations, with the possible exception of the results from the two Balkan countries.

**Figure 5: Social media source in the first education query**

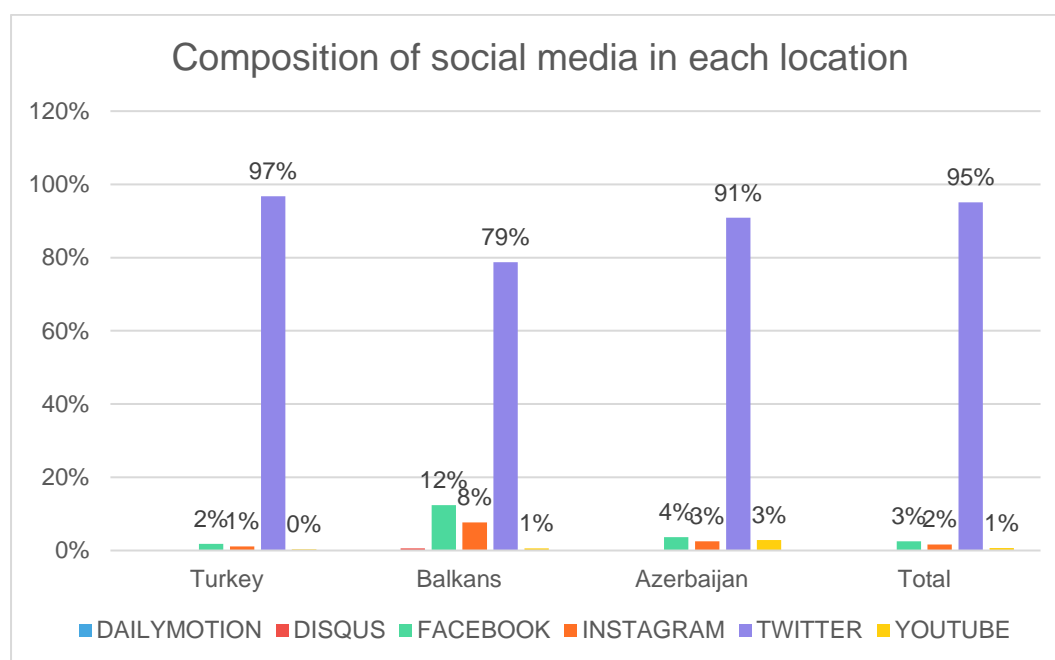
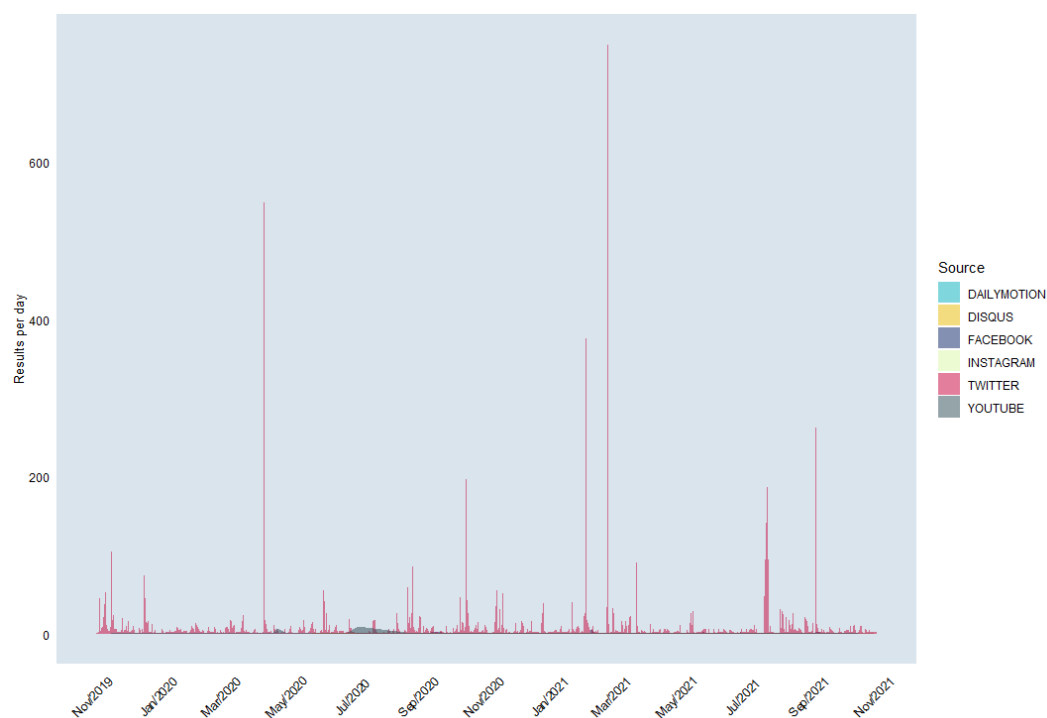


Figure 6 shows how results varied over time during the period we observed in the first query. There is a set of key dates when results spiked on Twitter, which is the main source of the data we examine here: March/April 2020, at the beginning of the global pandemic; October 2020; February 2021; July/August 2021; and September 2021. When looking at Turkey alone and the development of COVID-19 cases and daily deaths there, it becomes clear that the spikes we observe in the social media data roughly correspond to the pattern of the pandemic: in terms of daily cases, the first wave peaked in Turkey in April 2020; there was a second increase in cases and daily deaths between August and December 2020; a third wave between February and April 2021; and a further increase in cases and deaths between July and October 2021.<sup>11</sup>

<sup>11</sup> There has since been another wave in 2022, but this period is not covered by the analysis here. See <https://covid19.who.int/region/euro/country/tr> for a summary of how COVID-19 cases developed in Turkey over the course of 2020 and 2021.



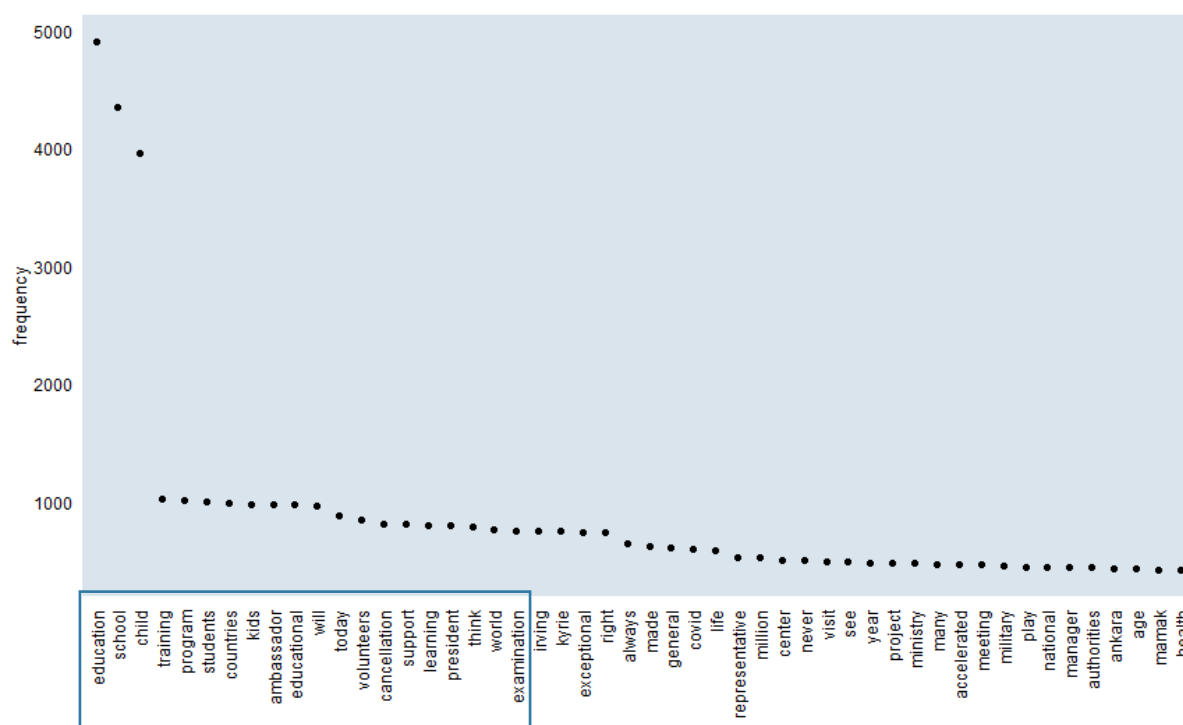
**Figure 6 Social media posts in the first education query over time**

### Content analysis: frequency of word occurrence, sentiment analysis, and topic modelling

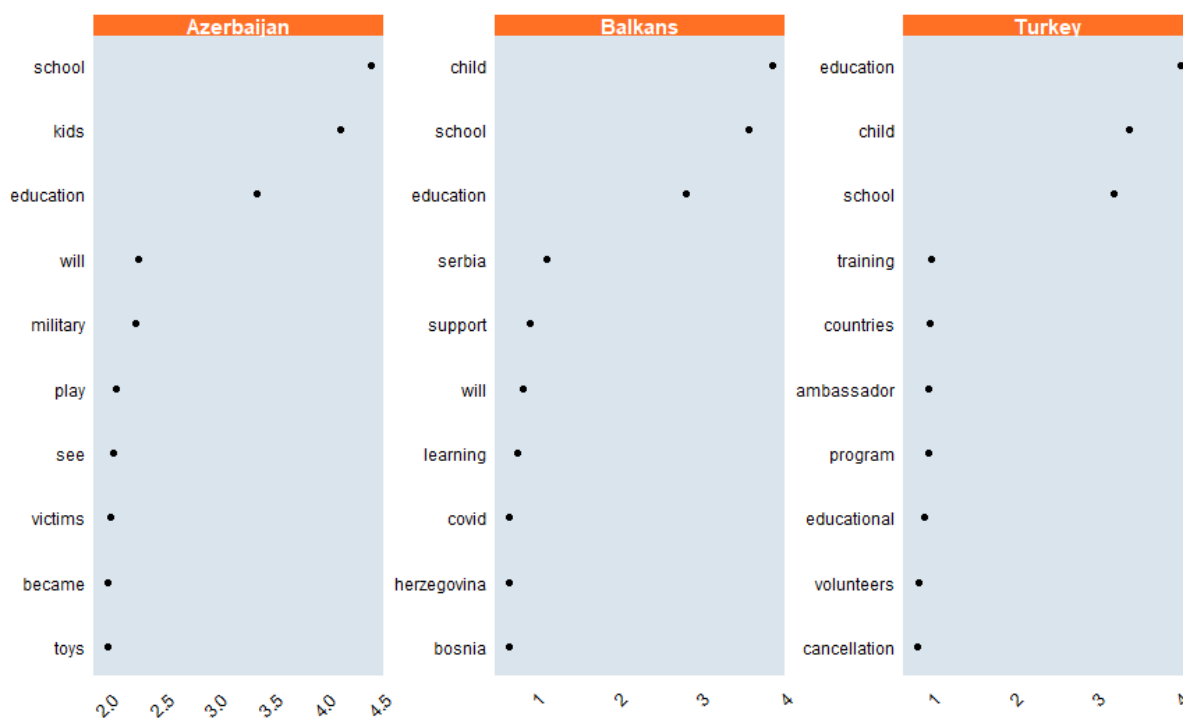
As mentioned in Section 2.4, a first analytical step was to look at the frequency with which certain words appeared in the text corpus. Figure 7 depicts the absolute frequency with which words appear in the data derived from our search query (after pre-processing described in Section 2.3). Given we are looking at an education query, these results are not particularly surprising at first sight: the most common words are 'education', 'school', and 'child'. However, among the 20 most common words, we also see terms relating to the disruption that COVID-19 created: 'cancellation' and 'support'.

We also look at these word frequencies by country and present the results in Figure 8. As before, we group Serbia and Bosnia & Herzegovina together under 'Balkans'. In this case, we plot the relative word frequency, i.e. how frequently the most frequent words appear in comparison to the total count of words per country group. We show the top 10 words per country group. Unsurprisingly, the top three words are, as before, 'education', 'school', and 'child' or 'kids'. However, there are some differences across these groups that already hint at how the conversation varies on the education topic: Azerbaijan is the only country where the 'military' and 'victims' play a role. In the two Balkan countries, 'support' is often mentioned, while in Turkey there is talk of 'programmes', 'volunteers', and 'cancellation'.

**Figure 7: Education—absolute wordcount for words with five or more mentions**



**Figure 8: Education—relative word frequency by country group**



To explore these word associations in more detail and in greater depth, we implement a topic modelling exercise (as described in Section 2.4) using LDA modelling. We develop a model that identifies a list of six topics, for which the top 10 terms are listed in Table 3. The results show that the differences between topics are sometimes subtle and not easy to discern. We have highlighted some words that we consider to be relevant for each topic.

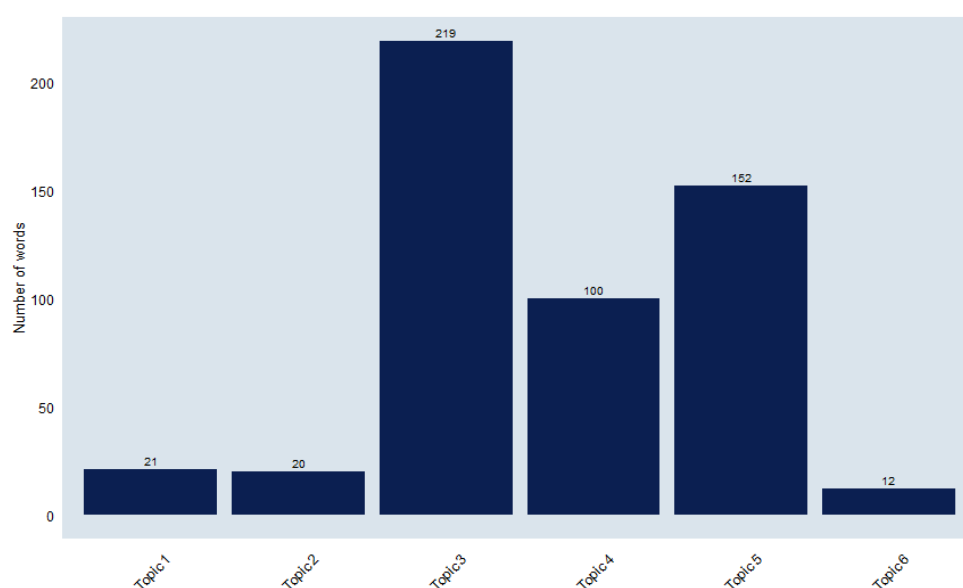
Topic 1 relates to military themes. Similarly, Topic 4 is about national and government programmes that support the right to school for students, with a possible discussion of the rights of refugees. Topics 3 and 5 have an international outlook, while Topic 3 is more about teachers and Topic 5 is about girls and the support they need. Topic 6 mentions ‘volunteers’, ‘ambassador’, and the ‘president’. Topic 2, finally, relates more to programmes that support children in their learning.

**Table 3: Education—topic modelling**

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
school	education	education	child	child	education
kids	<b>today</b>	school	education	school	<b>ambassador</b>
will	child	students	school	education	educational
never	<b>always</b>	child	<b>right</b>	<b>world</b>	<b>volunteers</b>
see	<b>learning</b>	<b>teachers</b>	<b>support</b>	<b>covid</b>	<b>countries</b>
<b>authorities</b>	visit	<b>international</b>	<b>program</b>	<b>million</b>	think
<b>victims</b>	<b>program</b>	<b>attack</b>	<b>ministry</b>	<b>girls</b>	<b>president</b>
<b>military</b>	<b>center</b>	<b>protect</b>	<b>national</b>	learning	examination
smile	accelerated	<b>training</b>	<b>financial</b>	every	kyrie
play	<b>representative</b>	support	refugee	face	irving

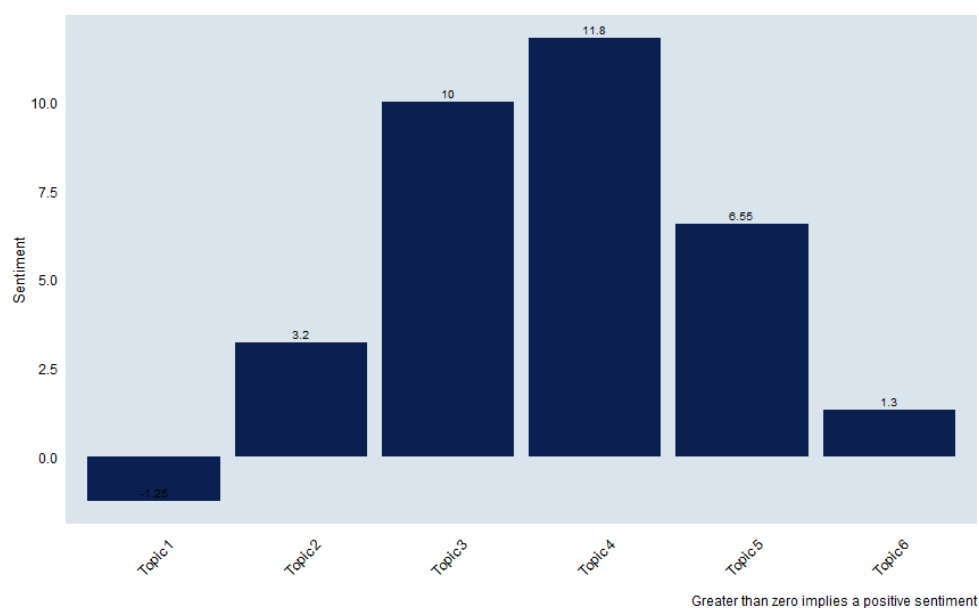
To explore these topics further, we look at additional analyses that shed light onto how these topics are made up and how they relate to the different countries from which we are drawing insights for the topic area of education. First, we look at how many words matter for each topic, i.e. how many words are needed to describe each topic. We can do this because topic modelling assigns ‘weight’ to words in terms of how much they explain each topic. In Figure 9, we graph the number of words that explain 60% of each topic. As can be seen, Topics 1, 2, and 6 need very few words to be explained.

**Figure 9: Education—number of words that explain 60% of each topic**



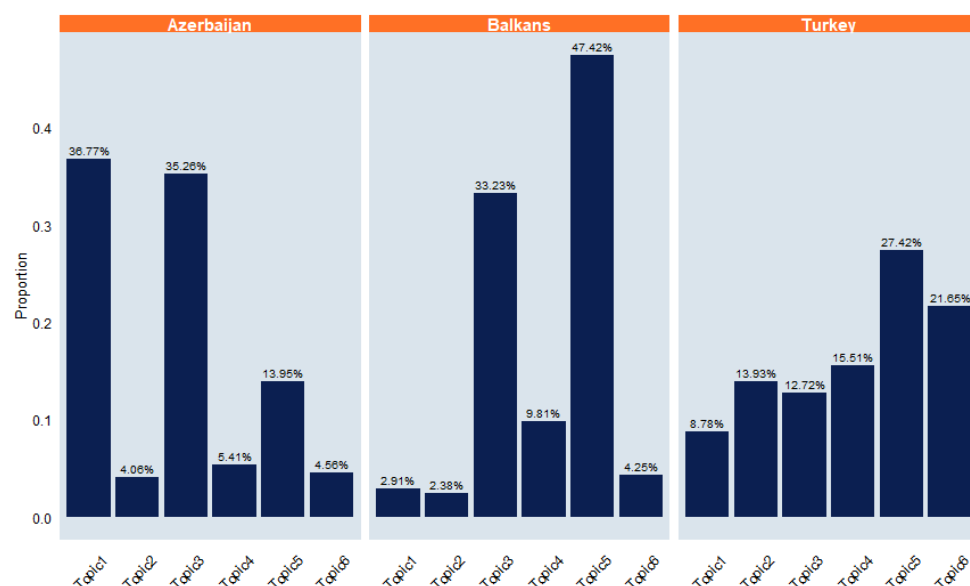
Second, we apply sentiment analysis to these topics to identify whether there are positive or negative sentiments associated with the words that primarily make up these topics. This analysis is, again, on the words that make up 60% of the explanatory value of each topic. The results are presented in Figure 10. Values that are larger than zero represent 'positive' sentiments. On average, we can see that only Topic 1 is associated with 'negative' sentiments. Topic 4 is overwhelmingly positive, which might not be surprising when looking at the list of words in Table 3, when it became clear that this topic relates to government support for schooling.

**Figure 10: Education—sentiment of topics**



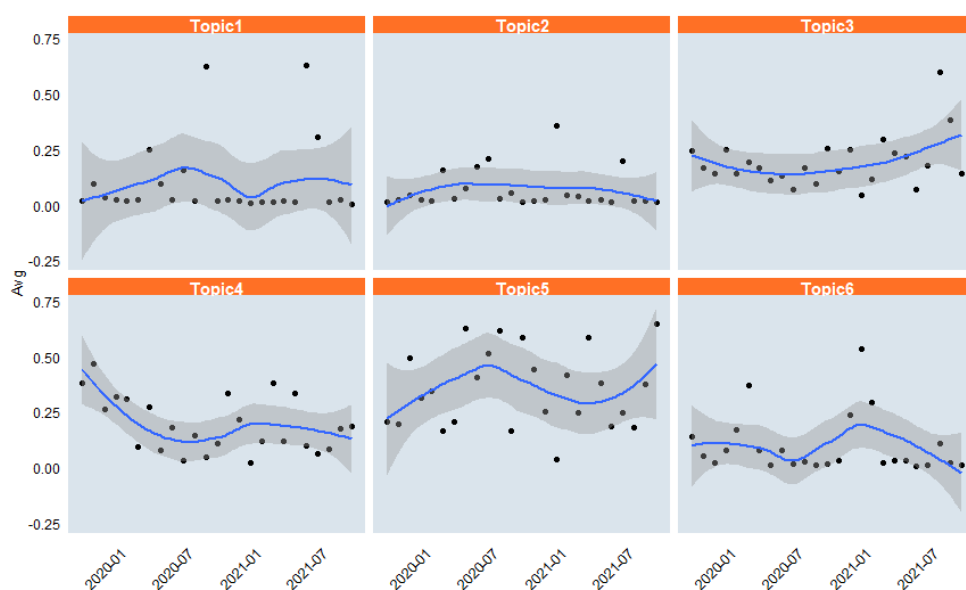
Third, we analysed how these topics appear across the three relevant country groupings for this query. In Figure 11, we can see that Topic 1 is only really of significant relevance in Azerbaijan, which—given the conflict in Nagorno-Karabakh—explains how this 'military' topic appears in the data. It also explains the negative sentiment identified above. In Azerbaijan, this conflict features heavily in social media discussions. In the Balkans, Topic 5 is most prevalent, relating to girls' education. Topic 3 appears relevant both in Azerbaijan and in the Balkans—an international outlook relating to teachers. Finally, in Turkey, Topics 5 and 6 feature significantly.

**Figure 11: Education—topics by country group**



Finally, we look at how topic relevance trended over the period we observed here. Figure 12 shows how prevalent the topics were in posts by month, together with a moving average (the blue line). The grey areas are confidence intervals around this average. The following key findings can be identified from these graphs. It is clear that the importance of Topic 4 decreased over time. This seems to imply that concerns over, or discussions of, financing and national support decreased over time as well. What started to matter more was Topic 5, which related to support for learning across the world, with a focus on girls’ needs. Similarly, Topic 3 increased in prevalence, which means that the needs of teachers were increasingly featuring in conversations we captured. Topics 1 and 6 increased in importance at some point but then decreased again over time, ending up pretty much where they started. The relevance of Topic 2 was relatively stable over time.

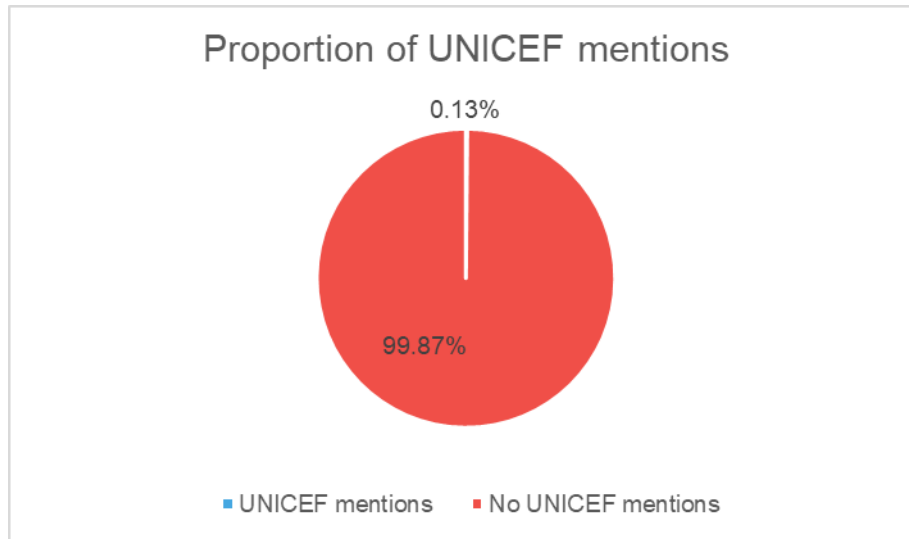
**Figure 12: Education—topic importance over time**



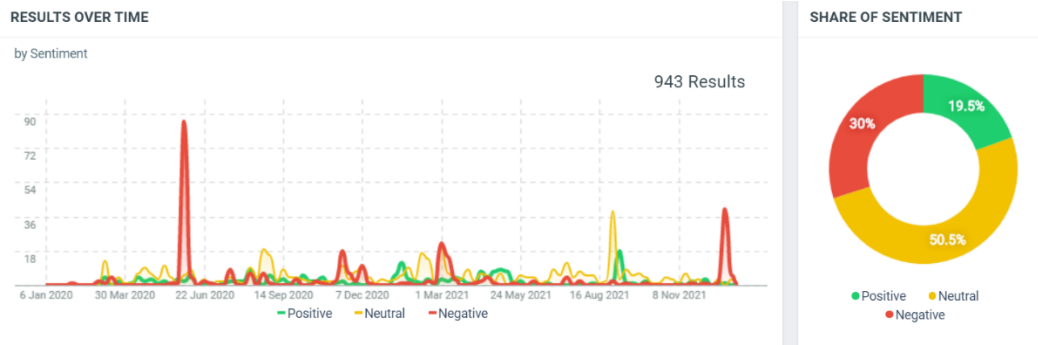




**Figure 15: Proportion of UNICEF mentions in non-UNICEF query on TW**



**Figure 16: Education—sentiment in UNICEF-related posts on the TW platform**





## 3.2 SP

The objectives of this chapter are to serve as a case study of the application of SML to understanding changes in the social media discourse in two countries: **Albania and Tajikistan**. In terms of sequencing, this analysis was implemented after the education analysis, presented in Section 3.1. The same analytical techniques were employed here, although results are presented in a slightly different way. In general, our analysis was guided by the following key questions, which are similar to the questions explored in the previous chapter:

- what are the key themes of SP discourse identified on social media?
- how has the SP discourse changed over the two years under investigation?
- what is the sentiment of this discourse? and
- what is the place of UNICEF in this discourse?

These questions were defined to contribute to the broader RTA questions, while accounting for the specific nature of social media data and the exploratory nature of the SML analysis. As mentioned in Section 2, our analysis employed the following approaches: frequency of mentions, frequency of themes, word clouds, networks/clustering of words, topic modelling, and sentiment analysis.

The overall period covered by the analysis was between 10 January 2020 and 26 December 2021. However, for the purposes of this specific SP analysis, we distinguish between the following subperiods:

- **Period 1—pre-COVID-19:** between 10 January 2020 and the first case of COVID-19 recorded in each country (08 March in Albania and 30 April in Tajikistan);
- **Period 2—early COVID-19:** between the first case of COVID-19 and 31 December 2020; and
- **Period 3—later COVID-19:** between 01 January 2021 and 26 December 2021.

In both countries, we were therefore able to distinguish between the pre-COVID-19 SP discourse, the early COVID-19 SP discourse, and the later COVID-19 SP discourse. We will refer back to these periods throughout our analysis.

### 3.2.1 The SP query

The first step of our analysis was to delimit what counted as SP-related discourse, or in other words to identify posts and other social media content that could be labelled as SP-related. For the SP SML work, we settled on the following definition: mentions of words in List A in proximity to words in List B, whereby List A refers to the 'social' dimension (words like 'social', 'humanitarian', 'poverty', 'vulnerability', etc.) and List B refers to the 'protection' dimension (words like 'protection', 'program', 'support', etc.). 'Proximity of occurrence' was defined as the occurrence of any word from List A within a two-word distance from any word in List B. Terms in Lists A and B were defined in English, Albanian, and Tajik. Albanian and

Tajik terms were verified and commented on by the SP specialists in the relevant UNICEF country offices.

This definition of SP-related social media content was translated into a search query using Boolean operators on TW, which allowed the platform to scrape social media and download content (we refer to individual bits of content as 'mentions') that met our definition. The query also specified the country from which the social media posts originated. Box 1 gives an example of such a query in the English language, which would apply to content written in English.

### Box 1: Example of a query that defines SP-related social media content on TW

```
(((social OR humanitarian OR poverty OR vulnerab* OR poor OR destitut* OR 'in need'
OR disadvantaged OR needy OR excluded OR exclusion)

NEAR/2

(protection* OR benefit* OR beneficiar* OR support* OR assistanc* OR transfer* OR
servic* OR payment* OR expenditure* OR program* OR scheme* OR pension* OR
allowanc* OR insuranc* OR aid OR initiative* OR intervention* OR project* OR relief OR
'financial support' OR 'financial assistance'))

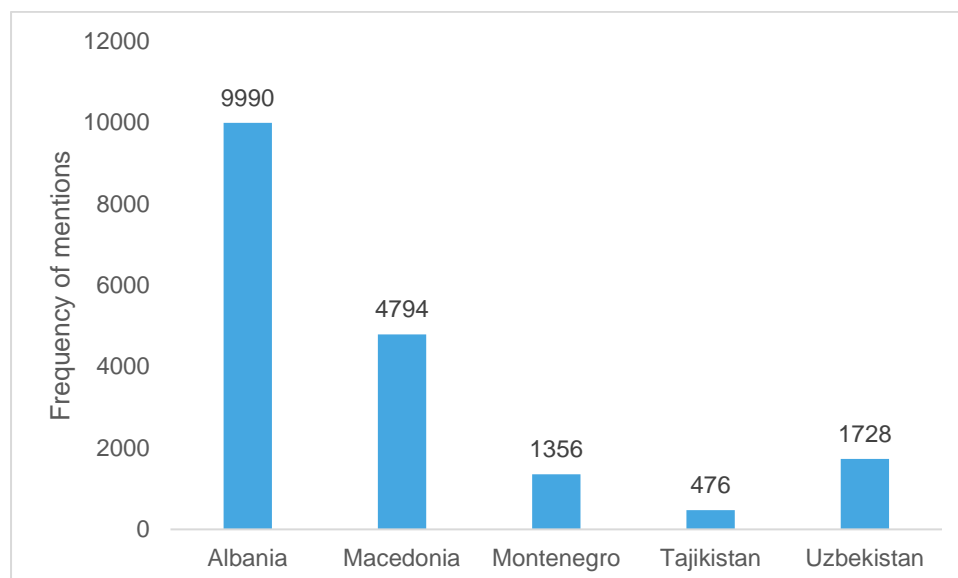
AND

(sourcecountry:mk OR sourcecountry:me OR sourcecountry:uz OR sourcecountry:tj OR
sourcecountry:al OR sourcegeo:mk OR sourcegeo:me OR sourcegeo:uz OR sourcegeo:tj
OR sourcegeo:al))
```

## 3.2.2 Describing the data

**Focusing Albania and Tajikistan gave us a chance to test our SML analytical approaches in very different social media and SP settings.** To put our case study countries into broader context, Figure 17 presents the number of posts mentioning SP in five ECARO countries that were subjects of the deep-dive RTA report on UNICEF's COVID-19 response in SP.<sup>12</sup> Albania and Tajikistan represent two extremes. In Albania, our query resulted in 9,990 mentions, while in Tajikistan, only 476 mentions were identified by the SP query over the nearly two-year period spanned by our study.

<sup>12</sup> Albania, Montenegro, North Macedonia, Tajikistan, and Uzbekistan.

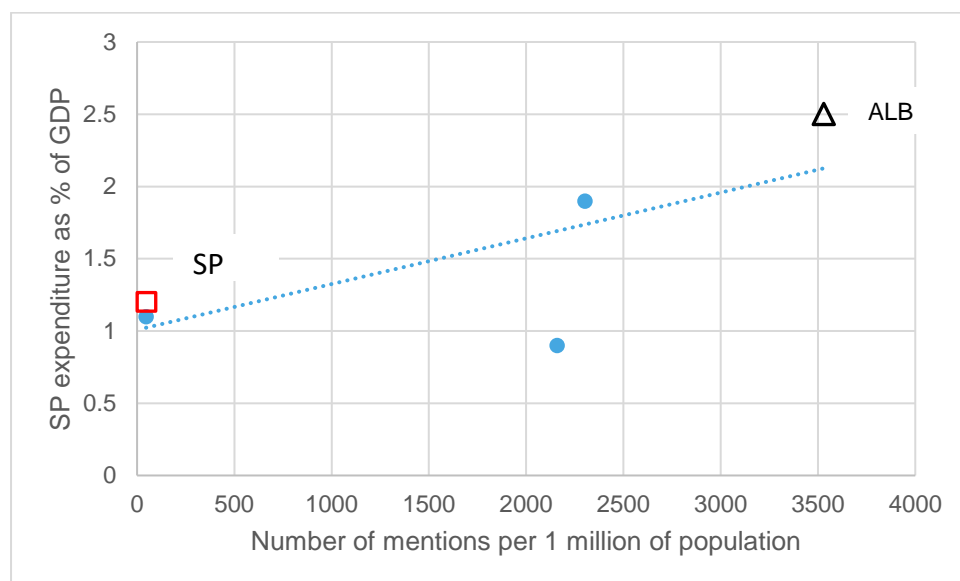
**Figure 17: Total number of occurrences of SP posts in select countries (pre-filtering)**

**The variation in the number of mentions is driven by many factors, but in our small sample it was correlated with a greater role of SP in the economies of the countries under discussion.** The frequency of SP mentions reflected a range of factors, from population size to the degree of participation in social media, control over social media content, the role played by the SP sector, and the urgency of SP-linked issues in the country. However, in the limited sample of five countries mentioned above, we observed a positive association between the frequency of SP mentions (adjusted for the size of the population) as the size of the SP sector<sup>13</sup> relative to the GDP<sup>14</sup> (Figure 18). In this regard, Tajikistan and Albania were at different ends of the spectrum but followed the trend: Tajikistan had a small SP sector and low number of SP posts, while Albania had a more active SP discourse and a relatively larger SP sector.

<sup>13</sup> The size of the SP sector is calculated as a sum of spending on social assistance and social pensions. Based on the ASPIRE database, the size of other programmes, such as school feeding and public works, in the five countries in question is negligible.

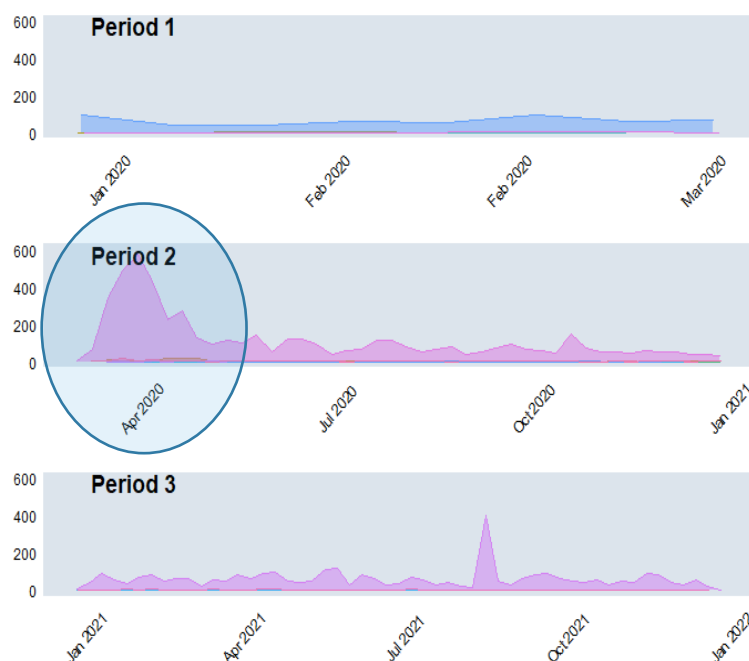
<sup>14</sup> Data from the ASPIRE database, a cross-country database of indicators of SP performance ([www.worldbank.org/en/data/datatopics/aspire/indicator/performance](http://www.worldbank.org/en/data/datatopics/aspire/indicator/performance)).

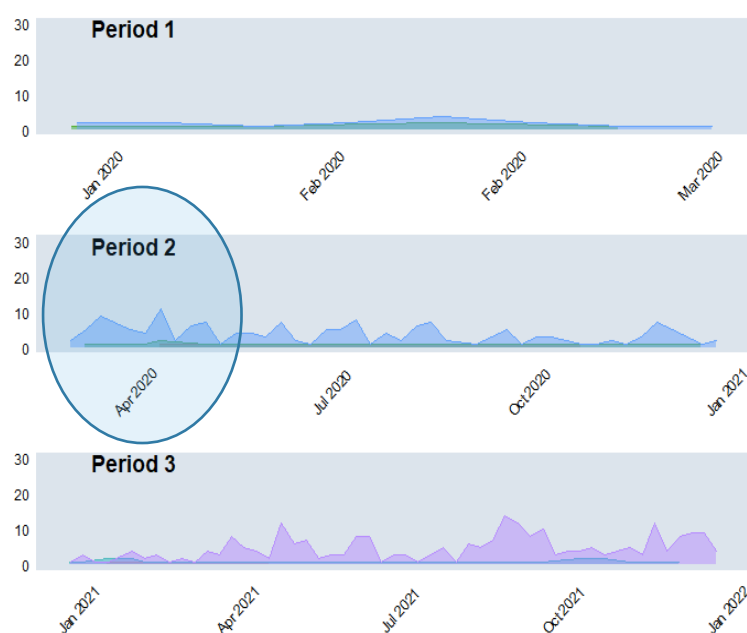
**Figure 18: Frequency of mentions of SP (population-adjusted) and relative size of SP sector (% of GDP)**



The frequency of SP mentions increased significantly both in Albania and in Tajikistan following the onset of the first wave of COVID-19. Figure 19 and Figure 20 clearly demonstrate a spike in SP discourse around the time the first case of COVID-19 was diagnosed in each country (08 March in Albania and 30 September in Tajikistan). In Albania, however, the upswing in the number of SP mentions was especially dramatic, but the spike subsided significantly over a two-month period from May onwards. In Tajikistan, neither the upswing nor the decline in the frequency of SP mentions was as pronounced.

**Figure 19: Evolutions in SP mentions in Albania by period**



**Figure 20: Evolution of SP mentions in Tajikistan by period**

### 3.2.3 Analysis results

#### Frequency of word mentions

The objective of this analysis is to describe how the SP discourse evolved over the three periods we identified by looking at the 20 most frequently mentioned words. In this and subsequent analysis, the data were pre-processed to remove stop words (articles, prepositions, conjunctions, personal pronouns, etc.). Our analysis looked at the normalised frequency of word mentions, which is achieved by dividing the number of mentions for a given word by the number of occurrences of the most common word, which tells us the prominence of words in the SP social media discourse.

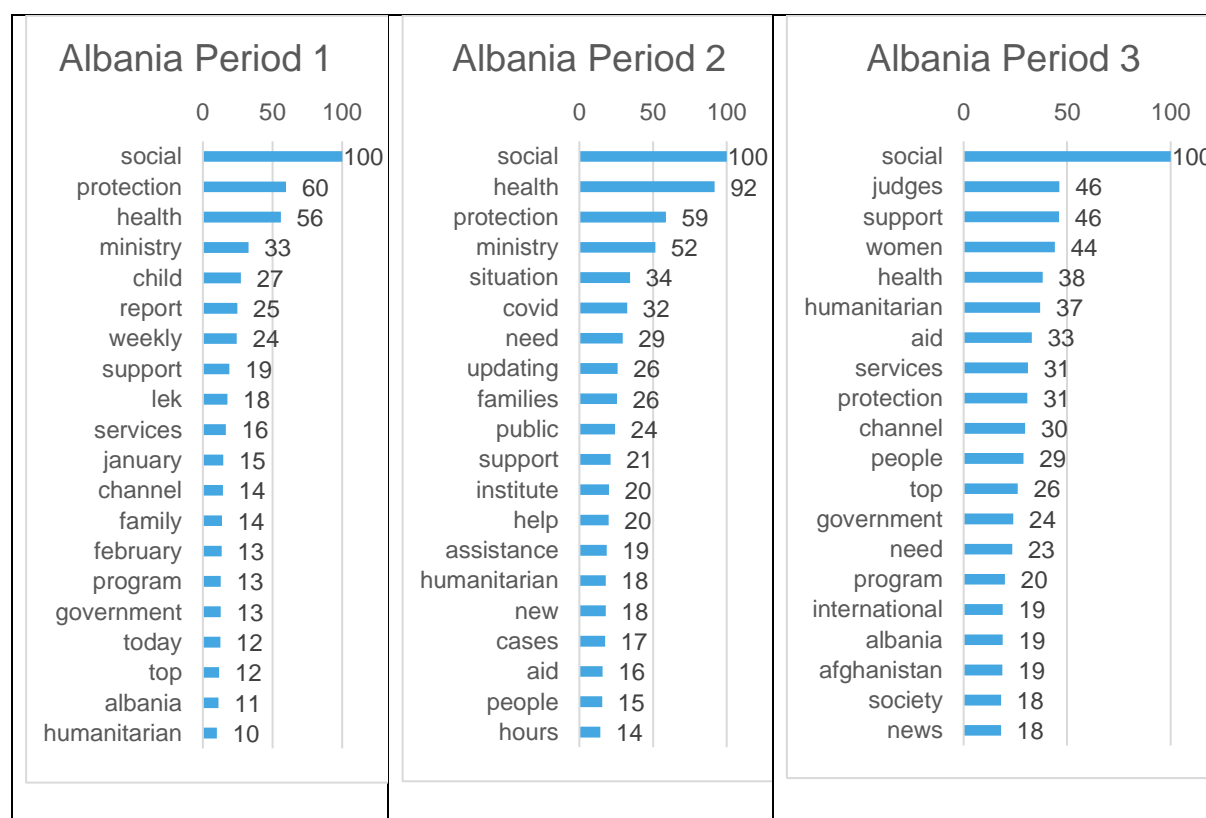
For example, an imaginary social media corpus for a given period has the word 'social' as the most common word with 200 occurrences, while the word 'child' has 50 occurrences in the same corpus. The normalised values of occurrences for 'social' would be  $f = 200/200 = 1$  and for 'child' it would be  $f = 50 / 200 = .25$ . Values 1 and .25 are then converted to a 0–100 scale, which enables us to say that the occurrence of the word 'child' in our example is equivalent to 25% of occurrences of the most common word in the corpus in a given period. This measure of prominence of words in the discourse is not affected by the total number of mentions in the discourse.

Word clouds represent another way to display most commonly occurring words for each period. They include the 70 most frequent words and reflect the importance of each word through the size and colour of their font. We used word clouds to illustrate the composition of the general SP discourse and the portion of that discourse that related to UNICEF (i.e. posts that mentioned either 'UNICEF' or 'The United Nations Children's Fund').

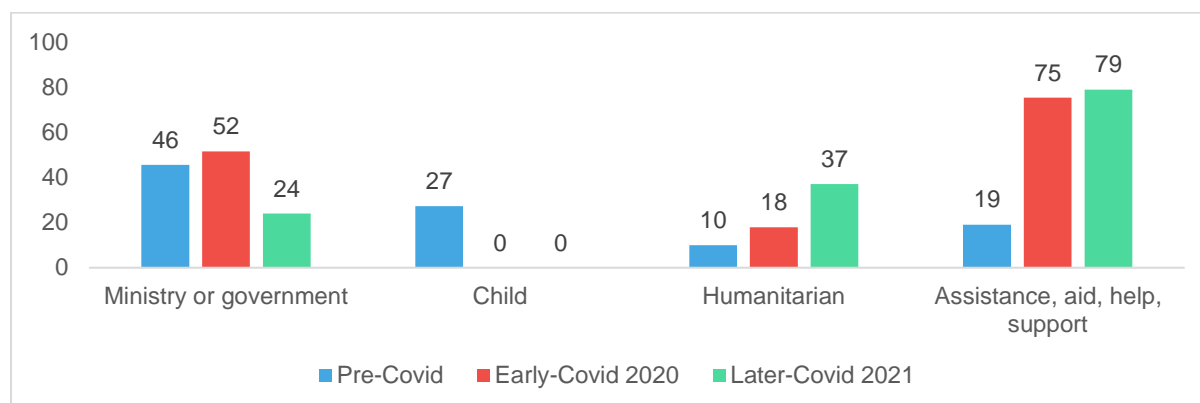
## Albania

In Albania, the general SP discourse expanded in terms of the volume of words and the range of issues that became central to it, with topics of assistance gaining prominence but with variable visibility of the role of government (Figure 21). Using the frequency output, we traced the centrality of specific words, when possible grouping them into cognate word clusters: e.g. 'ministry' and 'government' or 'assistance', 'aid', 'help', and 'support'. This analysis revealed that the relative frequency of mentions of 'ministry' or 'government' slightly increased during the early COVID-19 period but then declined significantly. In the meantime, mentions of humanitarian issues and of 'assistance', 'aid', 'help', and 'support' continued to climb, which potentially indicates that the importance of topics relating to assistance of all types was significant, but the role of the government and specifically the Ministry of SP and Health was becoming less visible in Albania over time as the COVID-19 pandemic evolved (Figure 22).

**Figure 21: The 20 most frequent words in Albania in each period (normalised frequency)**



**Figure 22: Albania—change in centrality of select word clusters (word frequencies normalised by occurrence of the most common word in the period)**



The monthly number of mentions of UNICEF remained constant before and after the arrival of COVID-19 to Albania (Table 4). Adjusting the number of mentions of UNICEF for the length of the period, we found no change in the frequency of mentioning UNICEF in the segment of social media we analysed. The rate remained steady at 3.3 to 3.6 mentions per month, according to the data we derived from our query. Given UNICEF’s increased involvement in the SP sector, this is surprising. The reason may be linked to the nature of the social media sources we were analysing: 90% of our Albania social media corpus consisted of Twitter messages, and it may well be the case that UNICEF’s visibility may not have expanded commensurately with its actual role in this less formal, less policy-focused or academic segment of social media. Whether this finding translates into any lessons for UNICEF depends on whether UNICEF views Twitter as a target audience of its information dissemination and visibility strategy.

**Table 4: Mentions of UNICEF in SP discourse**

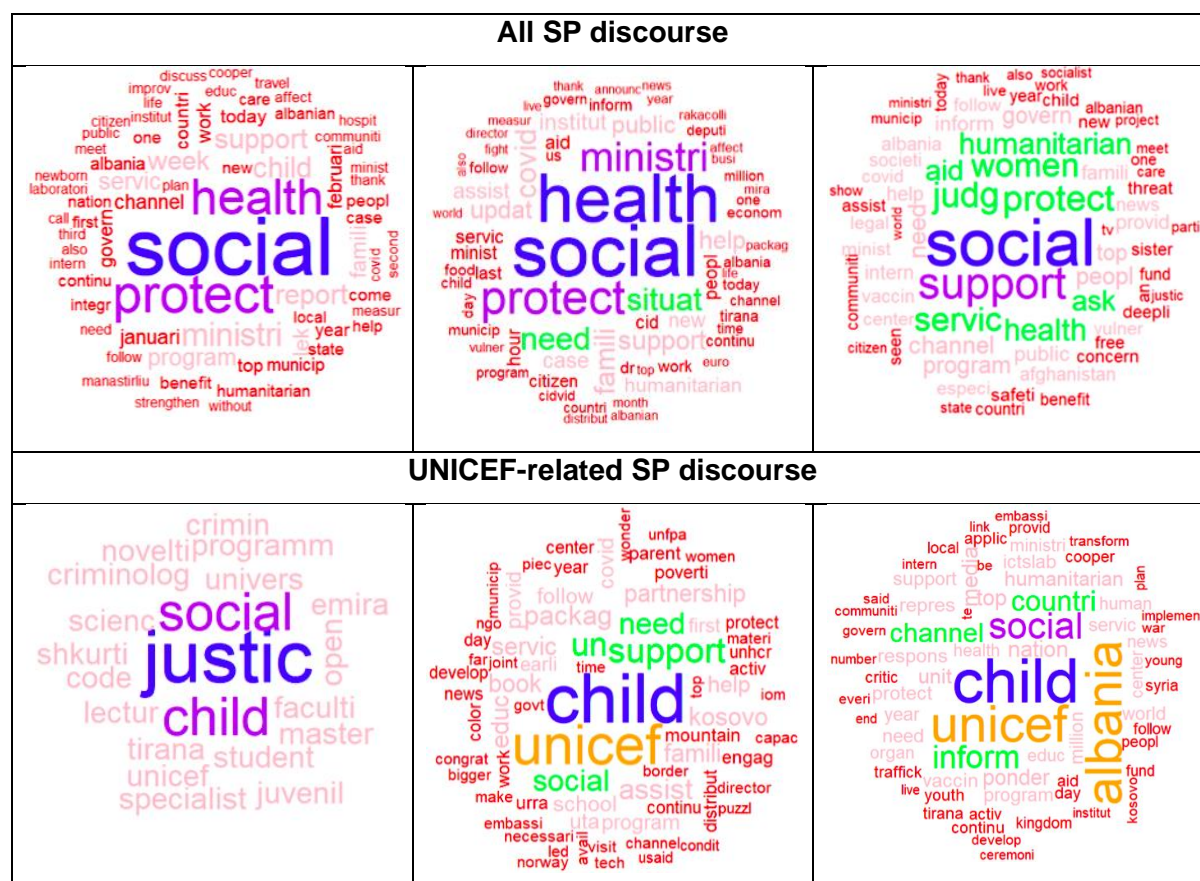
	Total number	Number per month
Period 1	7	3.6
Period 2	33	3.3
Period 3	39	3.3

In Albania, the UNICEF-related discourse over the two years under examination remained centred on issues relating to children, but shifted from issues relating to social support and assistance during the early pandemic period towards a focus on information dissemination and social media during the later pandemic period. We demonstrated these changes using the word cloud visualisation method. During the pre-pandemic months, the central word of the discourse mentioning UNICEF was ‘justice’ and a number of words related to education (‘university’, ‘master’s’, ‘student’), but the prevalence of ‘support’, ‘assistance’, ‘packages’,<sup>15</sup> and ‘services’ was visible during the early COVID-19 period. During the later COVID-19 period, words linked to information access—‘information’,

<sup>15</sup> For example, financial support packages.

'media', 'channel', 'program', 'news'—formed the most visible cluster. Looking at a side-by-side comparison with the general discourse, the UNICEF-related discourse seemed fairly well-aligned with the general SP discourse in its attention to various types of social support and assistance during the early pandemic period. However, during the later COVID-19 period, the general discourse still maintained a focus on social assistance, while UNICEF-related discourse largely lost that social support focus (even though the word 'humanitarian' did appear among the top 20 words).<sup>16</sup>

**Table 5: Albania—word clouds of top 70 terms central to all SP discourse and SP discourse mentioning UNICEF by period (from left to right, Periods 1, 2, and 3)**



To further explore the SP discourse in each of the three periods, we looked into the topic modelling results. Following the approach used in the education-related analysis and using the LDA approach, we identified six topics in each of the periods. As described in Section 2.4, topics—as identified by the topic modelling algorithm—are clusters of words that tend to occur together in the text. Interpretation of the list and naming or labelling the topic in a way a human would understand as such is a separate exercise. We interpreted and 'named the clusters' based on the top 10 words in each of the topics. Some words occurred across several topics and, in that case, we considered combinations of important words that set that topic apart from other adjacent or overlapping topics. The top 10 words for each topic for

<sup>16</sup> Please see Annex C for the word frequency graphs for Albania's UNICEF-related corpus.



each of the periods, and the interpretation/labelling of the topics based on them, is documented in Annex C.

**Over the two-year period covered by the study, four core themes characterised SP discourse in Albania: news and updates; humanitarian and social support and services; the situation in Kosovo; and international SP issues.** These four themes emerged over a total of 18 topics identified by the algorithm (six in each of the three periods) and showed the following evolution over time.

- News and updates over the three periods tracked issues in SP and health, initially (during the pre-COVID-19 period) focusing on the role of government in SP by discussing spending and issues related to children, but then switching to updates on the pandemic. During the third period, a specific person figured prominently in the updates: parliamentarian Ogerta Monastirliu.
- The discourse on humanitarian and social support and services evolved towards a greater attention to providing assistance to poor families, as well as towards the needs of, and support for, socially vulnerable groups such as children and girls specifically. The topics relating to the provision of services at the local/municipal level (including integrated social services provision), which were central during the pre-COVID-19 period, became less important during the later COVID-19 period.
- Kosovo-related SP issues were an important topic during the pre-COVID-19 and early COVID-19 periods, but lost their visibility during the later COVID-19 period.

**Table 6: Albania—topic modelling results**

	Pre-COVID-19 period	Early COVID-19 period	Later COVID-19 period
<b>News and updates on health and SP</b>	<p>1*. <i>Top Channel TV news on health and SP measures</i></p> <p>4. <i>Periodic reports of the Ministry of SP and Health touching on children and spending</i></p>	<p>5. <i>Top Channel TV news and information on the COVID-19 situation in Albania</i></p> <p>1. <i>Updates on the COVID-19 situation by the Ministry of SP and Health</i></p>	<p>3. <i>Top Channel TV news regarding health-related topic discussed by the Minister of Health and parliamentarian Ogerta Manastirliu</i></p>
<b>Humanitarian and social support to poor and services</b>	<p>2. <b>Humanitarian and social aid</b> to people in connection with <b>COVID-19</b></p>	<p>3. <b>Humanitarian</b> and social assistance to <b>families in need</b>, including economic support and food aid</p> <p>2. SP, services, and support for <b>children</b> and <b>vulnerable</b> households affected by COVID-19</p>	<p>1. <b>Humanitarian</b> aid and social support to <b>poor</b> families</p> <p>4. SP and health services and programmes for the <b>vulnerable</b></p>

	<p>3. Provision of integrated SP <b>services</b> related to <b>health</b> by <b>municipalities</b></p> <p>6. <i>Social support services in communities, including services to children and care</i></p>	<p>5. Social and legal aid for the <b>vulnerable</b>, including <b>girls</b></p> <p>6. Social support for <b>people in need</b>, including <b>children</b></p> <p>6. <i>Social and humanitarian aid provided by/to the <b>municipality</b> of <b>Tirana</b> related to <b>housing</b></i></p>
<b>Kosovo SP issues</b>	<p>5. International and government issues related to <b>Kosovo</b></p>	<p>4. Social and humanitarian assistance for the <b>poor</b> in <b>Kosovo</b></p>
<b>International SP issues</b>		<p>2. <i>Issues related to <b>justice</b> with respect to <b>women</b> in <b>Afghanistan</b></i></p>

\* Numbers 1–6 are assigned to the topic in each period by the topic modelling algorithm.

UNICEF was referenced in a meaningful way<sup>17</sup> during the earlier COVID-19 period in connection with the topics of SP, services, and support for children and vulnerable households affected by COVID-19. During the later COVID-19 period, it was referenced under the topic of Top Channel TV news regarding health-related topics discussed by the Minister of Health and parliamentarian Ogerta Manastirliu. In the former case, UNICEF appeared in the first 19% of most significant words and in the latter among the top 40% of most significant words for the respective topics. This indicates that UNICEF is more central to the former topic, but peripheral to the latter.

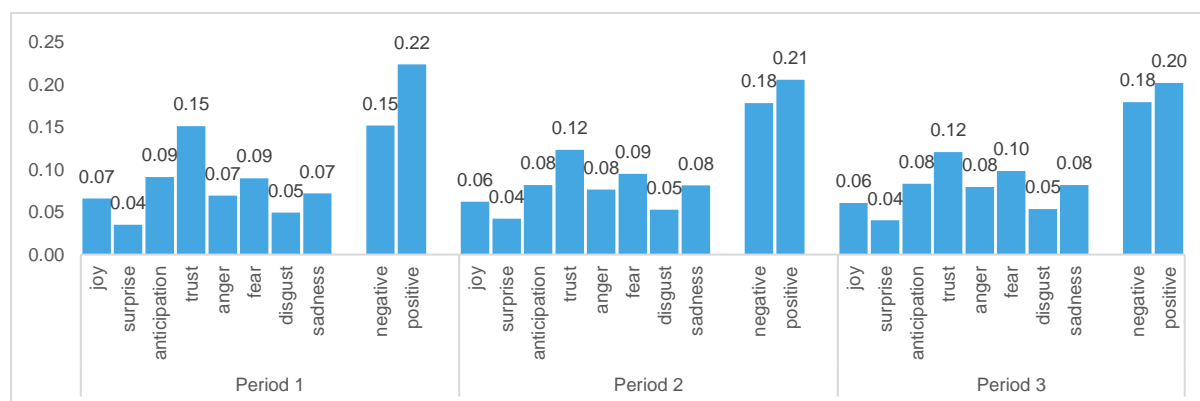
The sentiment analysis proceeded by assigning a sentiment score to each word in a text and then computing an aggregate score for each text (or post). The scoring factors assigned to each word were obtained from the National Research Council Canada (NRC) sentiment and emotion lexicons, which is commonly used in the NLP literature. The lexicon assigns polarity and intensity scores to each word in the English language. Polarity, or directionality, of the sentiment is conveyed by a sign (plus or minus), while its intensity is conveyed by the absolute value of the score. For instance, the word 'great' has a polarity of +1.2, while the word 'acceptable' has a polarity of +0.1, meaning that both have a positive polarity but the former is more intensely positive than the latter. In addition to scoring words on the positive–negative dimension, the NRC lexicon assigns words with respect to specific emotions: surprise, disgust, anger, sadness, joy, fear, anticipation, trust.

The sentiment profile of the discourse over the three periods is summarised in Figure 23. The profile was largely stable over the two years under investigation, with a few exceptions: the degree of trust declined following the pre-COVID-19 period. Positive sentiments still

<sup>17</sup> We define 'meaningful mention of UNICEF' under a topic if the words 'UNICEF' or 'United Nations Children's Fund' appear in the list of the top words accounting for 60% of the topic.

dominated negative sentiments in the SP discourse over the three periods, but the net positivity diminished over time.

**Figure 23: Sentiment and emotion profile of the Albania's SP discourse in three study periods (scoring based on the NRC lexicon)**

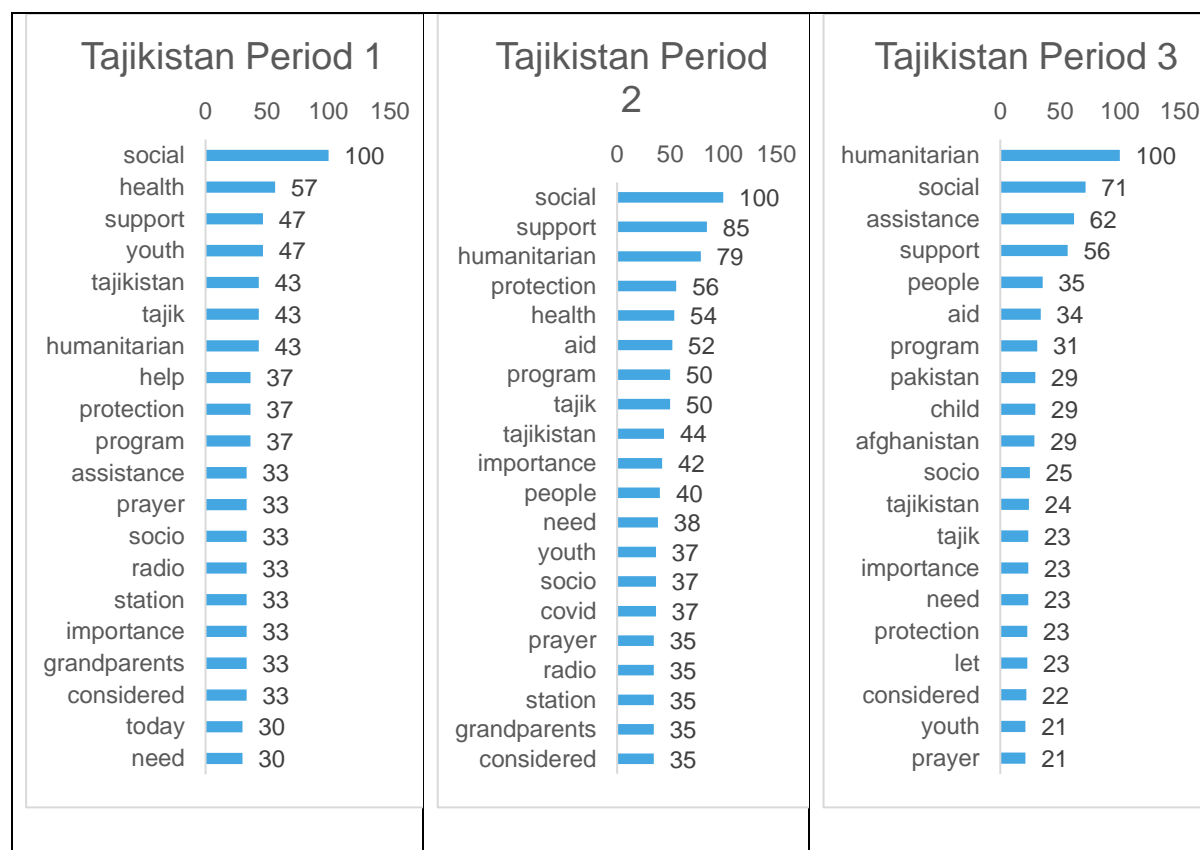


## Tajikistan

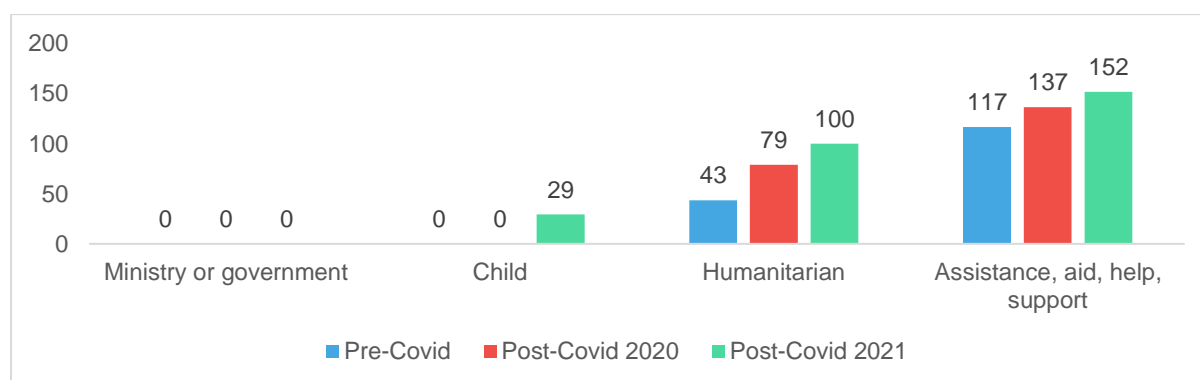
For Tajikistan, as can be seen in Figure 24, following the start of the pandemic, the SP discourse not only expanded in terms of the overall volume of comments but became more diverse, focusing on a wide range of key words. During the later COVID-19 period, the SP discourse crystallised around humanitarian and social assistance. For instance, prior to COVID-19, the most relevant seven words in the SP discourse had a relative frequency of 40% or more (i.e. they occurred at least 40% as often as the most frequent word). During the early COVID-19 period, this number increased to 11 words, with a relative frequency of 40% or more; but during the later COVID-19 period, there were only four such dominant words.

**In Tajikistan, the topics of humanitarian and social assistance were central in the general SP discourse prior to the arrival of the pandemic, but later gained even further prominence.** Analysis of frequencies showed that mentions of humanitarian issues and of 'assistance', 'aid', 'help', and 'support' steadily increased during the pandemic. However, mentions of the government or its ministries (e.g. the Ministry of SP and Health) did not appear prominently in this discourse. Mentions of children gained centrality during the later COVID-19 period (Figure 25).

**Figure 24: The 20 most frequent words in Tajikistan in each period (normalised frequency)**



**Figure 25: Tajikistan—change in centrality of select words (word frequencies normalised by occurrence of the most common word in the period)**



The discourse that mentioned UNICEF was very limited in volume during the pre-COVID-19 and early COVID-19 period. Only during the later COVID-19 period were the number of mentions of UNICEF sufficient to make meaningful inferences about their content (Figure 26). The UNICEF-related discourse during the later COVID-19 period focused on children, SP, support, and services, but also frequently mentioned migration. Thus, its focus on social assistance (broadly understood) was consistent with the general SP discourse. At the same time, the UNICEF-related discourse paid attention to the issue of migrants and children in a way the rest of the SP discourse did not.

**Figure 26: Tajikistan—word clouds of up to top 70 terms central to all SP discourse and SP discourse mentioning UNICEF**

Period 1	Period 2	Period 3
<i>General SP discourse</i>		
<i>UNICEF-related SP discourse</i>		
<p>Insufficient data on UNICEF</p>	<p>No mentions of UNICEF in text</p>	

Over the two-year period covered by the study, four core themes characterised the SP discourse in Tajikistan: traditional radio programming on youth-related issues; humanitarian assistance; the government’s SP response; and support for peace in connection with political processes in Afghanistan. These four themes were the result of groupings of more than 18 topics produced by the topic unsupervised modelling procedure for the three periods under study. Specific topics that gained prominence over time show how these themes evolved.

- Traditional radio programming on youth was a prominent topic running throughout the three periods. The topic touched on social issues and formally met our definition for SP-related content, but it seemed to take cultural and religious perspectives (the term ‘prayer’ figured prominently). Still, the fact that this topic was consistently present and showed considerable stability points to a robust cultural/religious discourse that may be worth considering when communicating on SP.
- The theme of humanitarian assistance across the three periods centred on references to humanitarian aid in connection with crisis response and support for vulnerable groups. The International Federation of the Red Cross and Red Crescent (IFRC) was mentioned during the early period, which reflects the IFRC’s contribution to crisis response. References to other countries in the region (Afghanistan, India, and Pakistan) were common, pointing to a regional perspective on the evolution of the pandemic and the COVID-19 response.

- The discourse on government response during the period prior to the first case of COVID-19 mentioned health supplies, the effect of COVID-19 on vulnerable populations, and social care provisions in Dushanbe. During the early COVID-19 period, two different topics under government response invoked 'cooperation': one referenced cooperation in connection with the protection of children, and the other mentioned cooperation in connection with financial support. During the later COVID-19 period, the government response centred around vulnerable groups and COVID-19 and the protection of children and migrants. References to multiple development partners in connection with the government response (IFRC, the United Nations Population Fund (UNFPA), the European Union, and France) further suggest the importance of cooperation to the SP discourse in Tajikistan.
- During the later COVID-19 period, the topic of social media coverage regarding support for peace efforts in Tajikistan in connection with political instability in Afghanistan (which coincided with the withdrawal of the United States army from the country and the Taliban takeover) became prominent.

**Table 7: Tajikistan—topic modelling results**

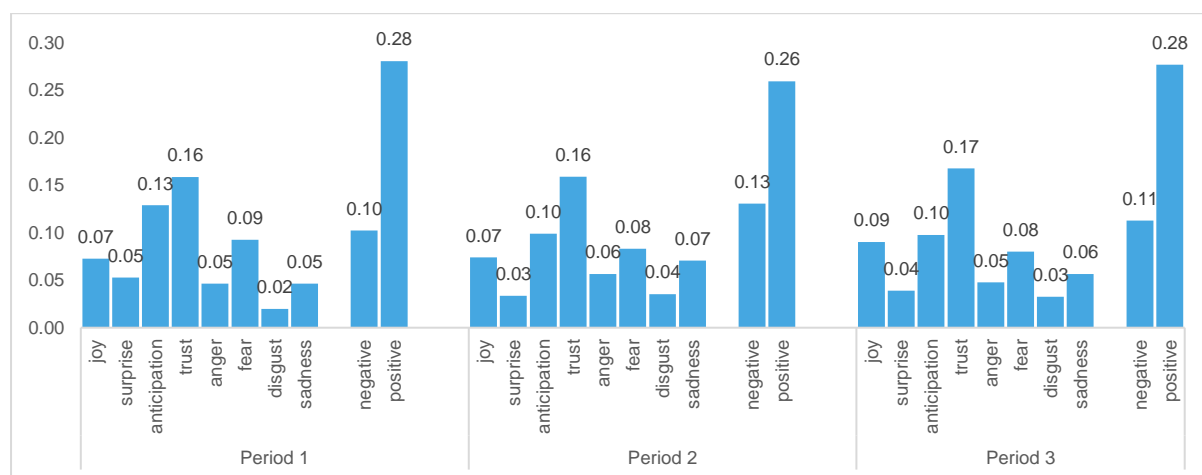
	Pre-COVID-19 period	Early COVID-19 period	Later COVID-19 period
<b>Traditional radio programming on youth</b>	2*. Traditional religious youth-centred radio programme	6. Traditional religious youth-centred radio programme	1. Traditional religious youth-centred radio programme
<b>Humanitarian assistance</b>	1. Humanitarian need and assistance/aid in <b>crisis</b>  3. Minister of Health and SP and <b>IFRC's</b> humanitarian role	1. Humanitarian assistance/aid for <b>vulnerable people</b> around the world  5. Humanitarian aid and <b>medical</b> support to fight against <b>COVID-19</b> in Tajikistan and <b>India</b>	2. Humanitarian assistance in the face of urgent need (mentioning <b>Afghanistan and Pakistan</b> )  4. Humanitarian assistance and support, including provision of food (mentioning <b>Pakistan</b> )
<b>Government's SP response</b>	4. Ministry's SP response and <b>health supplies</b>  5. Support to <b>vulnerable population</b> affected by <b>COVID-19</b> (mentions <b>France</b> )  6. Social care provision in <b>Dushanbe, facilities, training,</b> and attention to <b>psychic</b> health, (mentions <b>UNFPA</b> )	2. Ministry of SP and health: protection of <b>children</b> and <b>cooperation</b>  3. Government SP support and programmes linked to <b>COVID-19 impact</b>  4. <b>Financial</b> support to address health and social needs and ministry's <b>cooperation</b>	3. SP of <b>vulnerable groups</b> in connection with <b>COVID-19</b> (mentioning <b>European Union</b> and <b>Central Asian countries</b> )  5. SP and services provided by the ministry to <b>children and migrants</b>
<b>Support for peace</b>			6. <b>Social media</b> coverage of aid to <b>support peace</b> efforts in Tajikistan as linked with events in <b>Afghanistan</b>

\* Numbers 1–6 are assigned to the topic in each period by the topic modelling algorithm.

'UNICEF' appeared as an important word during the pre-COVID-19 period under the topic of support to the vulnerable population and again during the later COVID-19 period under the topic of SP and services provided by the ministry to children and migrants. In the former case, 'UNICEF' appeared in the first 50% of most significant words and, in the latter case, it appeared among the top 26% of most significant words for the respective topics.

The sentiment profile of the discourse over the three periods is summarised in Figure 27. The emotion profile was largely stable over the two years under investigation. The sentiment was overwhelmingly positive, although the *net* positive effect initially declined from 0.18 points to 0.13 points during the early COVID-19 period.

**Figure 27: Sentiment and emotion profile of Tajikistan's SP discourse in three study periods (scoring based on the NRC lexicon)**



## 4 Conclusions, lessons learned, and recommendations

As mentioned in Section 1, the objective of this workstream, as part of the larger RTA of UNICEF's COVID-19 response in the ECARO region, was to explore the use of NLP on social media data. We have presented results applied to two topic areas: education (Section 3.1) and SP (Section **Error! Reference source not found.**). In this section, we present conclusions for each of these topic areas, lessons learned from this work, and recommendations on the way forward for using SML and NLP for evaluation purposes.

### 4.1 Conclusions

#### 4.1.1 Education

- The data analysed for the purposes of this report were social media posts, predominantly sourced from Twitter and originating from Turkey. Together, these constituted about 80% of the overall text corpus included here.
- According to the results from Query 1, social media posts relating to education and COVID-19 in Turkey, Azerbaijan, Bosnia & Herzegovina, and Serbia expanded over time, with spikes that roughly corresponded to waves of the COVID-19 pandemic. Education as an issue relating to COVID-19 seems to have increased in importance over time.
- At an aggregate level and across the countries covered by the first education query, our word frequency analysis revealed that the conversation on education included discussions around training sessions, programmes, support, government officials, and cancellations, and involved an international outlook. This general perspective, however, hid differences across countries: in Azerbaijan, and in contrast to the Balkans and Turkey, words relating to military conflict were of importance.
- Our topic analysis reflected some of these country-level differences. We identified a topic, strongly associated with negative sentiments, that dealt with military activities and victims and that was really only of significance to the conversation in Azerbaijan. Clearly, the education conversation here was also influenced by the ongoing conflict.
- Other topics were less clearly differentiable. However, we identified two topics that seem to have increased in importance over time and that are of importance in all three of the country groups we investigated. On the one hand, these related to protecting, supporting, and training teachers. On the other, they related to how students, perhaps with a focus on girls, needed support for learning in the face of COVID-19. Both topics were positive in terms of their sentiment and they seemed to capture a large variety of issues relating to education and COVID-19 in general. Clearly, focusing on issues relating to these topics (e.g. in communications or other activities by UNICEF) would correspond to what matters for this conversation on social media in the countries we included in this analysis.



- Interestingly, in a broader query on education and COVID-19 in these same countries, UNICEF-related posts made up only a very small proportion of the overall text corpus. This seems to indicate that UNICEF does not necessarily feature very prominently in this conversation.

#### 4.1.2 SP

- Both in Albania and Tajikistan, the SP discourse expanded in terms of the overall volume of comments and became more diverse following the arrival of COVID-19 in the countries. The range of topics central to the SP discourse expanded, as evidenced by the greater number of frequently mentioned words.
- In **Albania**, the theme of social and humanitarian assistance consistently gained prominence over the study period, but the centrality or visibility of the role of the government during the later COVID-19 period declined somewhat. In **Tajikistan**, the SP discourse crystallised around humanitarian and social assistance, but the government was not frequently mentioned.
- In **Albania**, topic modelling showed that the core theme of social and humanitarian assistance, identifiable through frequency analysis, in fact involved multiple topics that evolved over time.
  - During the pandemic period, greater attention focused on giving assistance to poor families, as well as to the needs of and support for socially vulnerable groups, such as children and girls specifically. At the same time, topics relating to the provision of services at the local or municipal level (including integrated social services provision), which had been central during the pre-COVID-19 period, became less important during the later COVID-19 period. Further, Kosovo-related SP issues were an important topic during the pre-COVID-19 and early COVID-19 period, but lost their visibility during the later COVID-19 period.
- In **Tajikistan**, topic modelling showed the following.
  - The theme of humanitarian assistance across the three periods centred on support for vulnerable groups and had an important regional angle, with India, Pakistan, and Afghanistan being mentioned often.
  - The theme of government response, which was not easily identifiable through frequency analysis, was visible through topic modelling. The discourse on government response during the early COVID-19 period centred on cooperation in connection with the protection of children and cooperation in connection with financial support. During the later COVID-19 period, the government response centred around vulnerable groups and COVID-19 and the protection of children and migrants.
  - References to multiple development partners in connection with government response (IFRC, UNFPA, the European Union, and France) further pointed to the importance of cooperation in SP discourse in Tajikistan.

- The sentiment analysis showed that the sentiment was a net positive (positively coloured words outweighed the negatively ones). This positive net difference was especially significant in Tajikistan, perhaps because of the tendency in Tajikistan's social media to hold back criticism and dissatisfaction. However, even in Tajikistan, the positive net difference shrank during the early COVID-19 period but then recovered, likely due to the fact that the impact of COVID-19 was relatively minor in Tajikistan. In Albania, the positive net difference was not as great as in Tajikistan and shrank during the early and later COVID-19 periods. The emotion profiles showed few changes over time and were relatively similar across countries, with the exception of a more noticeable deterioration of trust in Albania.
- In Albania, the SP discourse made references to UNICEF.
  - UNICEF-related discourse over the two years under examination remained centred on issues relating to children, but shifted from issues relating to social support and assistance during the early pandemic period towards a focus on information dissemination and social media during the later pandemic period.
  - The UNICEF-related discourse seemed fairly well-aligned with the general SP discourse in its attention to various types of social support and assistance during the early pandemic period. However, during the later COVID-19 period, the general discourse maintained a focus on social assistance, while UNICEF-related discourse largely lost that social support focus.
  - The topic modelling showed that UNICEF was not referenced in a meaningful way in connection with key topics during the pre-COVID-19 period. During the earlier COVID-19 period, mentions of UNICEF appeared in connection with the topic of SP, services, and support for children and vulnerable households affected by COVID-19. During the later COVID-19 period, UNICEF was referenced most meaningfully under the topic of Top Channel TV news regarding health-related topics discussed by the Minister of Health and parliamentarian Ogerta Manastirliu.
- In Tajikistan, the SP discourse that referenced UNICEF:
  - was very limited in volume during the pre-COVID-19 and early COVID-19 periods; only during the later COVID-19 period were the number of mentions of UNICEF sufficient to make meaningful inferences about their content;
  - focused during the later COVID-19 period on children, SP, support, and services, but also frequently mentioned migration;
  - focused on social assistance (broadly understood) in a way that was consistent with the general SP discourse—at the same time, UNICEF-related discourse gave attention to the issue of migrants and children in a way the rest of the SP discourse did not; and
  - was most meaningful in the topic modelling analysis in relation to the topic of SP and services provided by the ministry to children and migrants, and was less meaningful in connection with the topic of providing support to the

vulnerable population affected by COVID-19. The association of UNICEF with other topics was largely tangential.

## 4.2 Lessons learned

In this section, we summarise some key lessons learned from the SML work implemented since August 2021. These lessons focus on the context within which this SML work was implemented, i.e. the TW platform used to retrieve data and the analytical work done with this data, but can also be applied to SML work more generally. We also focus on the usefulness of such analyses for the purposes of evaluations, in particular of UNICEF's work.

**Overall, one key lesson learned is that the role of NLP analysis of social media data as an evaluation tool still needs to be better defined.** NLP analysis of discourse is potentially a very promising area of research, but its value added as an instrument of evaluation still needs to be better explored. Clearly, if the object of evaluation is to explore and assess—as we have done—conversations relating to a specific topic area and whether UNICEF-related discourse is different from general discourse in certain areas (such as SP), it is feasible to gain insights that are relevant to this objective using NLP. For instance, in the SP topic area, our analysis identified meaningful differences between the two discourses using word frequency analysis and topic modelling.

However, if the object of evaluation is ultimately to make inferences about the effectiveness, efficiency, and relevance of UNICEF's response to the COVID-19 pandemic, then the NLP analysis would only be able to contribute part of the picture and it would be important to understand the limitations and place of this analysis. Further, additional analysis and triangulation would be required to build a link from these findings to meaningful and sound inferences about UNICEF's effectiveness, efficiency, and relevance. For instance, once we had identified that UNICEF's SP-related discourse in Albania did not centre much on the topics of social and humanitarian assistance, we would need to understand whether this was an indication of (for example) a lack of relevance or a strategic decision to allow other development partners or the government to take the lead in this discourse. Similarly, once we had identified that education-related conversations often did not mention UNICEF, we would need to investigate whether this indicated that UNICEF was focusing interventions on a specific area that is not of relevance in a country.

**A possible way forward could be to include SML in early stages of evaluations as a formative or exploratory tool informing the design and direction of an evaluation, rather than the other way round.** This is due to the exploratory nature of SML and the fact that topic modelling can help identify what matters in conversations online. Topics identified with this modelling—their trends over time and patterns across geographies (like the ones identified in this report)—could serve as a basis for defining evaluation questions and investigating the relevance of UNICEF's activities (or those of other institutions).

We faced a range of limitations from the TW platform, primarily relating to data content, accessibility, and the pre-processing required to conduct NLP. The proprietary algorithms we have referred to in Section 2 need to be mentioned here as well.

- **Limited content availability presented a problem for our analysis.** Around 70% to 80% of the posts were accessible in the form of snippets only, which usually consisted of

the first few sentences of a larger body of text, often because of firewall restrictions imposed by the platforms owning the content. The implications of this limited content are significant and need to be thought through. For instance, the initial lines of text may not be representative of the main topic of the text, to say nothing of the full range of topics that may be covered in the full article. As described above, for the purpose of this analysis, we chose to focus on full text content, which allowed us to be confident that our findings were not biased due to the nature of the snippets. At the same time, it allowed us to focus on a specific type of content only that was more 'conversational' and interactive in nature (most of the full text content came from Twitter simply because the content is not copyrighted and the posts tend to be shorter, well below the size of an average *New York Times* article), which may in fact give a focused picture of a real-time discourse. Still, going forward, it would be important to extend the analysis beyond snippets in order to capture a more comprehensive view of the online content.

- **The time horizon of the social media data available on the TW platform was limited.** Historical build-ups of more than two years come at an additional cost, which needs to be kept in mind in the future, when longer time-horizons are intended to be included in SML work based on TW data.
- **Data download restrictions in the form of per-day and per-download caps presented a difficulty.** For the present analysis, we solved this problem by manually downloading data on a daily basis in order to stay within the required limits. This process was laborious and required significant time investments.
- **The translation function for downloaded data in TW was not automated and required a separate step of uploading of native language texts to Google Sheets.** This additional step, again, required a manual process of ensuring data were properly fed into the Google Sheets set-up.
- **In-depth analysis and visualisation required significant programming and analytical inputs, including data pre-processing, which went beyond the capabilities of TW.** We found that the TW platform only performed minimal pre-processing, without removing stop words, stemming, etc., which meant that many of these steps had to be implemented separately using R code.
- **In sum, the limitations of the TW platform required the majority of the analysis to be done off-platform to enable the generation of insights relating to, for example, topic modelling.** A core methodology and the means of presenting the data off-platform using R were developed, but would benefit from additional fine-tuning and automation.

### 4.3 Methodological recommendations for further use of SML

In light of the above conclusions and lessons learned, we develop some key recommendations here. These recommendations focus on the further use of SML and NLP methods for the purposes of UNICEF ECARO's evaluation work, both in general and specifically with respect to the tools that were used in the present assignment.

- **First, we recommend further experimenting with and developing the methodology for using NLP and SML in relation to evaluation objectives.** The

results presented in this report indicate that some useful and interesting insights can be gained with respect to conversations online, the topics that matter to the online audience, how they vary by country and period, and the role played by UNICEF. The key question is, however, how these results can help achieve evaluation objectives, e.g. with respect to analysing the relevance, effectiveness, and efficiency of UNICEF's interventions and activities. We hypothesise that the exploratory value of SML could best be materialised as a formative research process at the beginning of evaluations, in other words to *inform* later evaluation activities by, for example, providing an indication of the types of question and topic area to be investigated. Similarly, SML could be used to evaluate specific interventions or activities that target the social media space, such as social media campaigns. In this context, questions of relevance and effectiveness could be answered more directly using online data.

- **Second, we recommend exploring topic modelling and sentiment analysis methods further, potentially investigating combinations with deep learning algorithms.** In our analysis, we used LDA topic modelling to identify topics that were of relevance for the conversations on education and SP in the countries that we covered with respect to each of these topic areas. This generated some useful insights, in particular when looking at how topic relevance changed over time or varied across countries. This could be explored further, for example by looking into algorithms like *lda2vec*.<sup>18</sup> These algorithms promise to provide a more nuanced insight into the topics prevalent in a particular text corpus. This might be particularly useful when applied to a richer set of data (as explained below). In a similar vein, sentiment analysis approaches are being improved regularly in NLP, and exploring these in more detail could yield more nuanced insights in the future.
- **Third, we recommend addressing the issue of 'text snippets' before embarking on further SML analyses.** As explained in several different parts of this report, our SML analysis was limited to few sources from which we could gain access to 'full text', mainly Twitter. In future, SML and NLP should be implemented on full text documents from a larger set of sources. Gaining access to such documents might be costlier than using the existing TW platform, but could generate insights that are more representative of the general online media landscape.
- **Fourth, we recommend also addressing the issues we encountered with respect to download restrictions and off-platform translations.** Where possible, this should be simplified and standardised in a way that does not require the manual processes implemented for our analysis.
- **Finally, once the above points have been addressed, we recommend automating some of the off-TW analysis implemented here.** In theory, it is possible to do this—e.g. via an API integration—but this would only be feasible once issues such as translation and download restrictions have been dealt with.

---

<sup>18</sup> See <https://medium.com/nanonets/topic-modeling-with-lsa-psla-lda-and-lda2vec-555ff65b0b05> for a brief explanation.

## Annex A: SML memo

### UNICEF ECARO COVID-19 RTA: SML

25 October 2021

---

#### Introduction

UNICEF's ECARO launched Phase 2 of RTA as a follow-up to Phase 1. In 2020, the UNICEF evaluation function launched and rolled out Phase 1 of the RTA (RTA-1) of the UNICEF response to COVID-19 at the country level. The objective of Phase 2 of the RTA is to inform a forward-looking reflection on the implementation of specific aspects of the country office response to COVID-19.

In Phase 1, OPM performed text analytics on the documents shared by UNICEF country offices that explained their response towards COVID-19. However, these documents presented a data point in time and, during Phase 2, it was decided by UNICEF and the research team to analyse real-time social media data around COVID-19 and UNICEF's country offices. A set of topics to be explored were selected, such as SP and education, for selected countries. This was referred to as SML.

SML analysis allows us to capture what is being discussed in social media on topics of interest. It is important to keep in mind that what is being said on social media is not equivalent to what is being done. We will be dealing with perception, visibility, and optics relating to SP and education. Favourable mentions of a programme do not necessarily indicate that the programme works well, but it may mean that it is perceived as working well.

Therefore, the RTA team looked for various SML tools that could provide us with access to real-time data. TW is one such SML tool. UNICEF suggested using TW as they already had a contract with them and were using TW on their other projects.

**At the beginning of this process, the expectation was that, with TW, the team could acquire the relevant SML data and use analytical tools to present those data in a useful pattern in order to tell a coherent story that could respond to the objectives of this work.**

#### Timeline of the work so far

- **Step 1—Introductory call with TW:** The TW team gave a demo of their platform to OPM and the UNICEF team on **11 August 2021**. They presented multiple utilities of their platform. The team was interested in conducting topic modelling and making word clouds and, based on their demo, assumed that the TW platform had all the required abilities. Jointly with UNICEF, the decision was taken to explore TW capabilities further and build up an analytical workplan using TW data.

- **Step 2—OPM team reviewed TW platform:** After the introductory call, TW account access was granted and the OPM team explored the TW platforms. This exploration mainly entailed preparing search queries in order to fetch most relevant data to answer evaluation questions. Furthermore, we explored various TW visualisation features.

We used TW documentations and the digital excellency centre<sup>19</sup> to familiarise ourselves with the platform.

- **Step 3—Email exchange and first follow-up meeting with TW:** After initial exploration, OPM exchanged emails with TW (**23 September 2021**) on a set of questions/clarifications. As a result, we set up our first follow-up discussion meeting with TW on **29 September 2021**.
- **Step 4—Second follow-up meeting with TW:** The first meeting with TW answered a few questions. The OPM team went back with TW solutions and realised that further clarifications were required. OPM therefore requested a second follow-up meeting with the technical team of TW. The meeting happened on **07 October 2021**, but no representative from TW's technical team joined this meeting.
- **Step 5—Email exchanges with TW:** The second meeting clarified significant limitations for OPM regarding the analytical rigour required for this project. Following this meeting, the OPM team continued having email exchanges with TW and is currently trying to set up a third meeting with the participation of the tech team. This is the step we are currently in.

## Findings and issues identified to date

In general, the OPM team has found that TW is an excellent resource when it comes to acquiring data from different social media resources. Their query writing tools are easy to use and interpret.

For example, in order to explore discussions on social media in Turkey about education and UNICEF in the context of COVID-19, the following query could be used:

```
(UNICEF AND (education OR student OR school OR "class room"~ OR Eğitim OR Öğrenci OR okul OR sınıf OR "sınıf odası"~ OR ((kovid OR covid OR pandemi) AND (Eğitim OR Öğrenci OR okul OR sınıf OR "sınıf odası"~)))) AND (sourcecountry:tr))
```

However, while exploring the results that this query (and others) yielded, the OPM team did come across a set of issues that would impact our analysis. We list these issues, and potential ways forward, below.

### Issue 1: Limited access to historical data

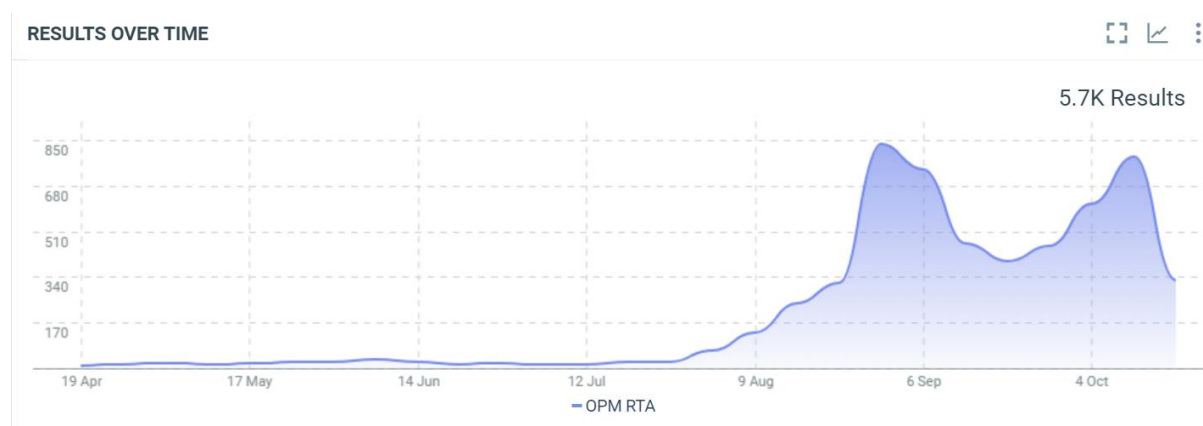
As mentioned above, one key objective for this SML work was to identify how conversations or discussions online about UNICEF, COVID-19, and related topics (e.g. such as education

<sup>19</sup> <https://talkwalker.digitalexcellencecenter.com/>.

and SP) changed over time during the period of the COVID-19 pandemic. This means that we need access to historical social media data (from the beginning of 2020) for our analyses.

We assumed that the results we were getting from TW were providing all the results inside the selected period. For example, if selecting a period of one year, we assumed that TW would provide data about that query for the last year. Our interest is to observe the period from the start of the pandemic up to the present. The OPM team therefore asked the TW team about this functionality during the first demo meeting (Step 1 above) and we were told it was feasible.

However, we quickly found this was not the case, at least in the standard query set-up, and TW only provides data for the last 30–60 days. As we wanted to analyse COVID-19 trends, this was a challenge.



**Solution to Issue 1:** TW informed us that there is a 'Historical Data' build-up solution that can provide historical data for up to five years, but this was a TW add-on and UNICEF would have to purchase it, which meant adding to their current contract.

Subsequently, our team was informed by TW that, with the current UNICEF contract, a historical build-up of two years was already available. Since then, OPM has repeatedly tried to build this data up but has not yet been successful. Currently, the issue is that it appears that the monthly quota for historical build-ups has been exhausted (erroneously it appears) by another project. This issue has yet to be solved.

## Issue 2: Pre-processing of data

As mentioned, the OPM team has been trying to come up with a coherent story around the perception on social media around COVID-19, UNICEF country offices, and the selected thematic areas (SP and education). To do this, we looked closely at the results from TW word clouds and other outputs. However, we realised that much of the outputs are clouded by irrelevant terms such as stop words ('a', 'is', 'the', 'them', etc.), and even the terms we used to search for results (e.g. 'UNICEF', 'education', etc.). This is a significant issue as we could not use the results from TW directly.

Just like in Phase 1 of the RTA, the OPM team was expecting to implement a thorough cleaning and processing of this data which includes splitting hyphenated words into two separate words; removing all numbers, punctuation marks, symbols (such as '\$', '@', '%'),

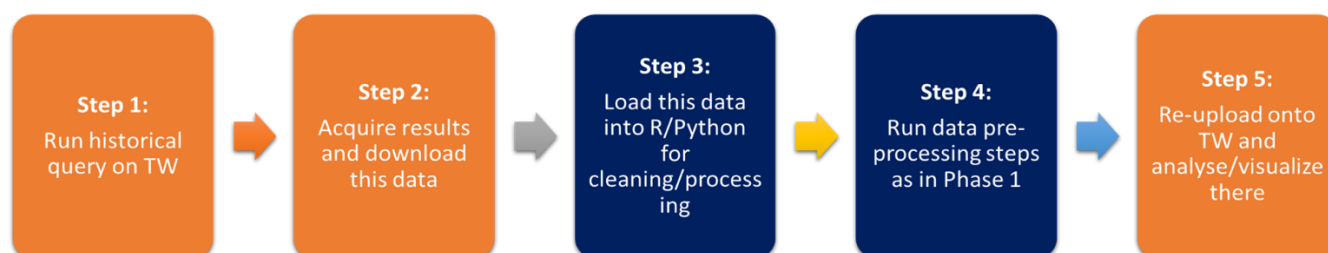


and URL websites; eliminating all 'stop words' from the body of text; and reducing words to their stem.

We had meetings with TW to discuss this pre-processing and cleaning of data. They suggested using more filters and manually cleaning the data. However, this is very laborious and not replicable, limiting our analytical capabilities.

**Solution to Issue 2:** OPM therefore proposed the following workflow to overcome this challenge (Figure 28). The processes contained in orange colour tiles would be run on the TW platform, while the ones in blue tiles would be run locally by OPM. We presented this process workflow TW and validated it with them.

**Figure 28: OPM's proposed process workflow to conduct data analysis**



### Issue 3: Uploading data back to TW

Following email exchange with TW, TW suggested that, to make user uploaded data visible on the TW platform, we should 'tag' it before uploading. This means adding another column in the uploaded data file. However, even after OPM followed the steps recommended by TW, the custom data uploaded by OPM seemed to 'disappear' on the TW platform.

**Solution to Issue 3:** We are currently in an email conversation with TW on this, but a solution to the issue has not yet been found.

### Issue 4: Downloaded data are not complete

In addition to the above problem, another issue with the workflow depicted in Figure 28 is that the data download ability of TW is currently limited in the sense that it does not allow downloading full text content from social media. **This means that Step 2 (downloading data) in the workflow of Figure 28 does not work as intended.** This was flagged as a caveat by TW tech team in an email on 12 October 2021.

The issue is that, if the result of a search is a long article, the downloaded data will only show what TW calls a 'snippet'. For example, a news article by a Turkish media outlet on education will only appear as few lines from the first paragraph of the article. See below for an example.

**Complete Article**

bianet'te önceki gün yer alan yazımda kurşun içeren boyaların çocuklarda kurşun maruziyetine yol açan en önemli kaynaklardan biri olduğunu belirtmişim. Bu mesele hakkında yeni bir yazı hazırlarken, okurlardan da çeşitli sorular geldi. Gelen sorulara da yanıt vermeye çalışarak çocuk sağlığı için önemli olan bu meseleye devam edeceğim.

**TIKLAYIN- Çocuklarda kurşun maruziyetine yol açan boyalar yasaklansın**

**Dünya genelinde durum nedir?**

Geçtiğimiz yıl UNICEF ve Pure Earth tarafından yayınlanan bir rapora göre dünya genelinde yaklaşık her üç çocuktan biri (yaklaşık 800 milyon) kurşun zehirlenmesine maruz kalıyor. Yapılan ölçümlerde çocukların kanındaki kurşun seviyesinin desilitre (litrenin onda biri, 100 mililitre) başına en az 5 mikrogram (µg/dL) olduğu belirlenmiş. Bu seviye, Dünya Sağlık Örgütü ve ABD Hastalık Kontrol ve Önleme Merkezi'nin çocuklarda kurşun maruziyetini azaltmak için bölgesel ve küresel ölçekte koruyucu-önleyici faaliyetlere geçilmesi için bir kriter olarak kabul ediliyor.

Hazırlanan raporda kurşunun çocukların beyinlerinde onarılamaz hasara neden olan güçlü bir nörotoksin olduğu belirtiliyor. Kurşun özellikle bebekler ve 5 yaş altı çocuklar için ömür boyu kalıcı nörolojik, bilişsel ve fiziksel bozulmalara yol açıyor. Çocukluk çaında kurşuna maruz kalma, zihinsel

**Downloaded Data Snippet**

bianet'te önceki gün yer alan yazımda kurşun içeren boyaların çocuklarda kurşun maruziyetine yol açan en önemli kaynaklardan biri olduğunu belirtmişim. Bu mesele hakkında yeni bir yazı hazırlarken, okurlardan da çeşitli sorular geldi. Gelen sorulara...

**This is a big issue**, as it makes our strategy of pre-processing data locally unfeasible (Step 4 in Figure 28). If we only download snippet data and clean it, we will not get the full picture of the conversations and issues we are analysing.

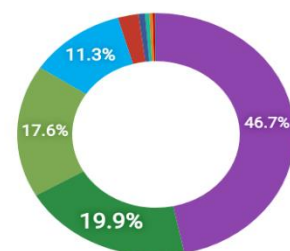
In order to understand more about this issue, the picture on the right states the media type from OPM search results. Of all search results on SP and education, 11.3% came from Twitter. We assume that most of the other media types are long posts that will definitely be affected by this snipped issue. We are not sure if Twitter (given the short number of characters in tweets) is affected but, overall, this graph shows that this issue is significant in any case.

The OPM team identified a set of key questions as a result of this. We put the following questions to TW and are awaiting their response:

- From which media types are we able to download complete or snippet data?
- Is there a word limit after which the results are exported as snippets?
- More importantly, is there a workaround where we can export full results using TW without manually copying and pasting from their source web links?

**Solution to Issue 4:** We emailed these questions to TW but still have to hear back.

#### SHARE OF MEDIA TYPES



## Possible ways forward

Given these issues, and the fact that some significant problems (regarding uploading data and the snippet issue) have not yet been resolved, the OPM team has been discussing what possible ways forward there could be.

**Option 1:** We only use TW's existing functionality, including historical data. If we go down this route, we will need to shift our analytical focus onto questions that do not require post-search processing of content. For instance, we can focus on:

- an analysis of the incidence of mentions of specific programmes or the incidence of other specific phenomena (**FEASIBLE**);
- an analysis of a change in the incidence of mentions over time (**LIKELY FEASIBLE**); and
- possibly a sentiment analysis, provided we are given more information on the ways in which sentiment analysis is done (**POSSIBLY NOT FEASIBLE**).

No topic modelling is feasible under this option due to TW's limited post-search data processing capability.

*This option is least work-intensive and is most likely to produce some results, but the level of research insights it will generate will be more limited. If we go down this route, we will need to adjust our research questions, which we have so far presumed will involve some form of topic modelling.*

**Option 2:** We use R to analyse exported data. The feasibility of this option depends on TW's ability to export meaningful amounts of text as opposed to only the 'snippets' it can currently export. TW's snippets do not have enough content to enable a meaningful topic modelling analysis. We would also need to work on visualisation tools.

If the data exporting issue is resolved, topic modelling will be feasible and will generate meaningful insight, but will involve considerably more work.

**Option 3:** Abandon the TW platform completely. At this stage, this will be the most labour-intensive and cost-intensive option. We do not consider this option in our way forward as it would be associated with a cost extension.

We propose the following for now:

- to manage our expectations about the usefulness of the results for our research; and
- to pursue both Option 1 and Option 2 as follows: in connection with Option 1, reframe our analysis purely in terms of analysing the incidence of terms, including the analysis of incidence overtime; in connection with Option 2, get clarity on the possibility of exporting larger snippets or full content and, if that proves possible, pursue Option 1. We aim to get clarity on the exporting issue within the coming week.

## Annex B: Search queries

- **Query 1: Education + COVID + UNICEF:**

```
(UNICEF AND (education OR student OR school OR 'class room'~ OR Eđitim OR Öđrenci OR okul OR sınıf OR 'sınıf odası'~ OR ((kovid OR covid* OR pandemi OR corona* OR korona*) AND (Eđitim OR Öđrenci OR okul OR sınıf OR 'sınıf odası'~))) AND (sourcecountry:tr)) OR
```

```
((UNICEF OR ЮНИСЕФ OR УНИЦЕФ) AND (education OR student OR school OR 'class room'~ OR obrazovanje OR školski OR škola OR učionica OR образование OR ученик OR школа OR школьный OR класс OR образовање OR школски OR учионица OR класна OR соба OR ((covid* OR ковид* OR pandemic* OR pandemija* OR пандемия* OR пандемија* OR corona* OR корона*) AND (obrazovanje OR školski OR škola OR učionica OR образование OR ученик OR школа OR школьный OR класс OR образовање OR школски OR учионица OR класна OR соба))) AND (sourcecountry:ba OR sourcecountry:rs)) OR
```

```
(UNICEF AND (education OR student OR school OR 'class room'~ OR təhsil OR tələbə OR məktəb OR 'sınıf otađı'~ OR ((covid* OR pandemic* OR pandemiya* OR corona* OR tac*) AND (təhsil OR tələbə OR məktəb OR 'sınıf otađı'~))) AND (sourcecountry:az))
```

- **Query 2: Education + COVID**

```
((education OR student OR school OR 'class room'~ OR Eđitim OR Öđrenci OR okul OR sınıf OR 'sınıf odası'~) AND (kovid OR covid* OR pandemi OR corona* OR korona*) AND (sourcecountry:tr) AND (sourcetype:SOCIALMEDIA_TWITTER OR sourcetype:SOCIALMEDIA_FACEBOOK OR sourcetype:SOCIALMEDIA_YOUTUBE OR sourcetype:SOCIALMEDIA_LINKEDIN OR sourcetype:SOCIALMEDIA_INSTAGRAM OR sourcetype:SOCIALMEDIA_MIXCLOUD OR sourcetype:SOCIALMEDIA_SOUNDCLLOUD OR sourcetype:SOCIALMEDIA_VIMEO OR sourcetype:SOCIALMEDIA_DAILYMOTION OR sourcetype:SOCIALMEDIA_VKONTAKTE OR sourcetype:SOCIALMEDIA_TWITCH OR sourcetype:SOCIALMEDIA_DISQUS)) OR
```

```
((education OR student OR school OR 'class room'~ OR obrazovanje OR školski OR škola OR učionica OR образование OR ученик OR школа OR школьный OR класс OR образовање OR школски OR учионица OR класна OR соба) AND (covid* OR ковид* OR pandemic* OR pandemija* OR пандемия* OR пандемија* OR corona* OR корона*) AND (sourcecountry:ba OR sourcecountry:rs) AND (sourcetype:SOCIALMEDIA_TWITTER OR sourcetype:SOCIALMEDIA_FACEBOOK OR sourcetype:SOCIALMEDIA_YOUTUBE OR sourcetype:SOCIALMEDIA_LINKEDIN OR sourcetype:SOCIALMEDIA_INSTAGRAM OR sourcetype:SOCIALMEDIA_MIXCLOUD OR sourcetype:SOCIALMEDIA_SOUNDCLLOUD OR sourcetype:SOCIALMEDIA_VIMEO OR sourcetype:SOCIALMEDIA_DAILYMOTION OR sourcetype:SOCIALMEDIA_VKONTAKTE OR sourcetype:SOCIALMEDIA_TWITCH OR sourcetype:SOCIALMEDIA_DISQUS)) OR
```

```
((education OR student OR school OR 'class room'~ OR təhsil OR tələbə OR məktəb OR 'sinif otağı'~) AND (covid* OR pandemic* OR pandemiya* OR corona* OR tac*) AND (sourcecountry:az) AND (sourcetype:SOCIALMEDIA_TWITTER OR sourcetype:SOCIALMEDIA_FACEBOOK OR sourcetype:SOCIALMEDIA_YOUTUBE OR sourcetype:SOCIALMEDIA_LINKEDIN OR sourcetype:SOCIALMEDIA_INSTAGRAM OR sourcetype:SOCIALMEDIA_MIXCLOUD OR sourcetype:SOCIALMEDIA_SOUNDCLLOUD OR sourcetype:SOCIALMEDIA_VIMEO OR sourcetype:SOCIALMEDIA_DAILYMOTION OR sourcetype:SOCIALMEDIA_VKONTAKTE OR sourcetype:SOCIALMEDIA_TWITCH OR sourcetype:SOCIALMEDIA_DISQUS))
```

- SP:

```
((social* OR humanitar* OR varfër* OR varfr* OR cənüşməria OR vulnerabël OR 'familje në nevojë' OR 'familjet në nevojë' OR 'personat në nevojë' OR 'njerëzit në nevojë' OR 'familje ne nevoje' OR 'familjet ne nevoje' OR 'personat ne nevoje' OR 'njerezit ne nevoje' OR përjasht* OR pambrojtur OR pafavoriz*) NEAR/2 (mbrojtj* OR përfit* OR mbështet* OR ndihm* OR transfer* OR shërbimet OR shërbimet OR pages* OR kesh OR program* OR pension* OR sigurimi OR shpenzimet OR pafavoriz* OR iniciativ* OR lehtësim OR projekt* OR subvencio* OR (ndihm* NEAR/1 financiar*) OR (paket* NEAR/1 financiar*)) AND (sourcecountry:al OR sourcegeo:al)) OR
```

```
((социјал* OR хуманитарн* OR сиромаш* OR рањив* OR искључен* OR socijaln* OR humanitarn* OR siromaš* OR ranjiv* OR isključen*) NEAR/2 (заштит* OR корис* OR подршк* OR помоћ* OR трансфер* OR услуг* OR плаћа* OR програм* OR пензи* OR осигура* OR расход* OR иницијатив* OR олакша* OR пројек* OR субвенци* OR (финансијск* NEAR/1 подршк*) OR (финансијск* NEAR/1 помоћ*) OR заштит* OR korist* OR podršk* OR pomoć* OR transfer* OR uslug* OR plaća* OR program* OR penzi* OR osigura* OR rashod* OR inicijativ* OR olakša* OR projek* OR subvenci* OR (finansijsk* NEAR/1 podršk*) OR (finansijsk* NEAR/1 pomoć*)) AND ( sourcecountry:me OR sourcegeo:me)) OR
```

```
((социјал* OR хуманитарн* OR сиромаш* OR ранлив* OR (социјал* NEAR/1 загрозе*) OR искључен*) NEAR/2 (заштит* OR корис* OR подршк* OR помош* OR трансфер* OR услуг* OR плаќа* OR програм* OR пензи* OR осигурув* OR расход* OR иницијатив* OR олеснув* OR проект* OR субвенци* OR (финансиск* NEAR/1 подршк*) OR (финансиск* NEAR/1 помош*)) AND (sourcecountry:mk OR sourcegeo:mk)) OR
```

```
((ичтим* OR ичтим* OR социал* OR башардўстона OR камбизоат* OR осебпазир* OR бечора OR осебпазир OR бенаво OR ниёзманд OR камбизоат OR ташаббус OR сабуќи OR лоиҳа OR истисно) NEAR/2 (хифз* OR ҳимоя OR фойд* OR бенефитсиар* OR дастгир* OR ёрдамчи OR интиқол* OR интиқолх* OR хизма* OR пардохт* OR барном* OR нафақ* OR сугурт* OR харочот OR субсиди* OR (дастгир* NEAR/1 молияв*) OR (ёри* NEAR/1 молияв*)) AND (sourcecountry:tj OR sourcegeo:tj)) OR
```

```
((ijtimoiy* OR gumanitar OR qashshoqlik OR zaiflik OR 'kambag'al' OR zaif OR muhtoj) NEAR/2 (himoy* OR foyd* OR 'qo'llab-quvvatlash' OR
```

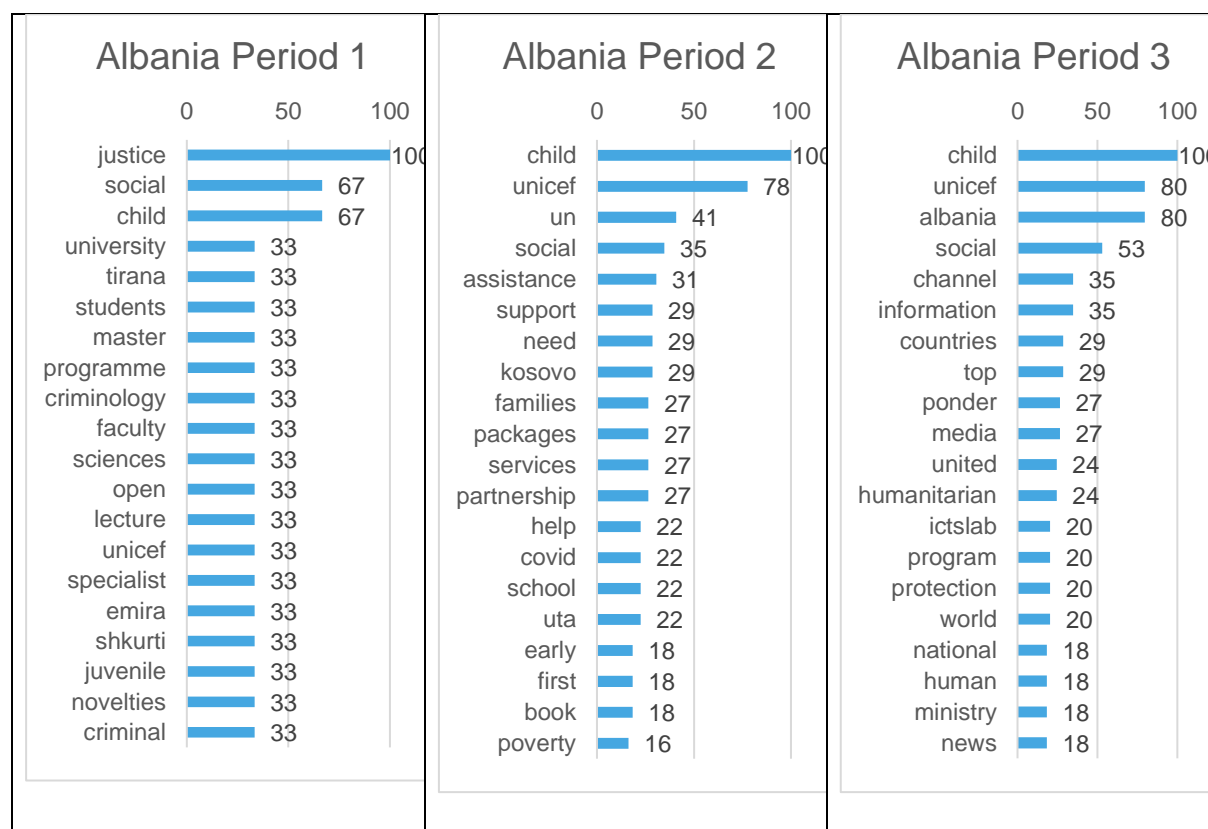
yordam OR 'o'tkazmoq' OR transfert\* OR xizma\* OR 'to'lov' OR dastur OR  
pensiya OR 'sug'urta' OR xaraja\* OR tashabbus OR loyiha OR nafaqa OR  
'moddiy yordam' OR 'kam ta'minlangan')) AND (sourcecountry:uz OR  
sourcegeo:uz)) OR

((social OR humanitarian OR poverty OR vulnerab\* OR poor OR destitut\*  
OR 'in need' OR disadvantaged OR needy OR excluded OR exclusion) NEAR/2  
(protection\* OR benefit\* OR beneficiar\* OR support\* OR assistanc\* OR  
transfer\* OR servic\* OR payment\* OR expenditure\* OR program\* OR scheme\*  
OR pension\* OR allowanc\* OR insuranc\* OR aid OR initiative\* OR  
intervention\* OR project\* OR relief OR 'financial support' OR 'financial  
assistance')) AND (sourcecountry:mk OR sourcecountry:me OR  
sourcecountry:uz OR sourcecountry:tj OR sourcecountry:al OR sourcegeo:mk  
OR sourcegeo:me OR sourcegeo:uz OR sourcegeo:tj OR sourcegeo:al))

## Annex C: Topic modelling for Albania and Tajikistan

### Albania frequency distribution of top 20 words

Figure 29: Normalised frequency distributions of top 20 words in the UNICEF-related discourse in Albania



### Topic modelling detail for Albania

To further explore the SP discourse in each of the three periods, we look to topic modelling. Using the LDA approach, we identify six topics in each of the periods. Topic modelling produced lists of words ranked in terms of their importance to the topic. These lists are clusters/patterns of words that tend to occur together in the text—this is what the algorithm identifies as a topic. Interpretation of the list and naming or labelling the topic a human would understand as such is an additional exercise.

To implement this labelling, we look at the top 10 words in each of the topics and identify the combinations of important words that set that topic apart from others, despite some overlap. We present the results of this topic modelling exercise below. The figures present, in detail, the words that constitute each topic, the proportion of overall text documents (comments)

that each topic constitutes in our text data, and the number of words that constitute 60% of the overall topic.

Period 1 topics can be defined as follows:

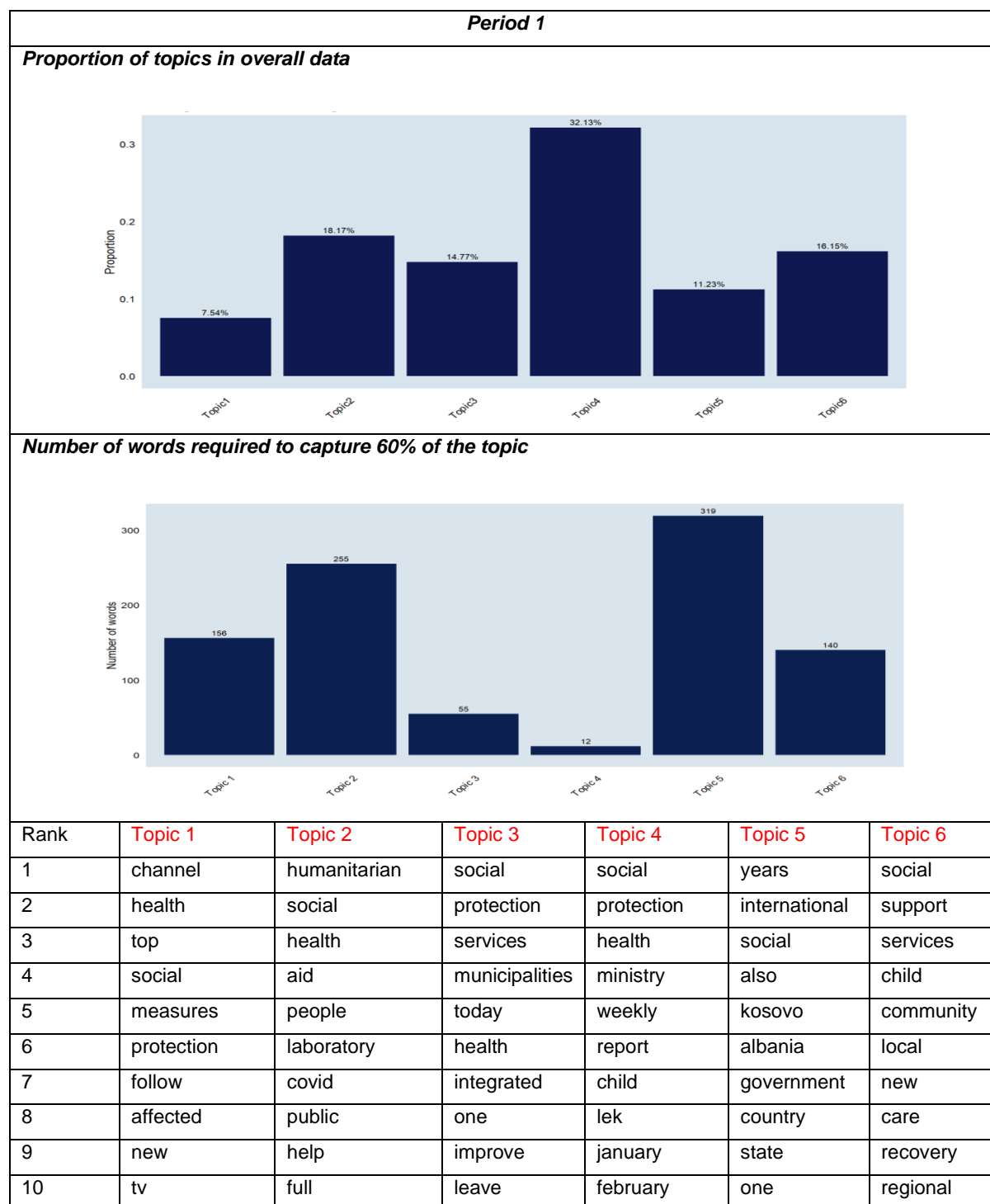
- (a) Top Channel<sup>20</sup> TV news on health and SP measures;
- (b) the provision of humanitarian and social aid to people in connection with COVID-19;
- (c) the provision of integrated SP services related to health by municipalities;
- (d) periodic reports from the Ministry of SP and Health touching on children and spending;
- (e) international and government issues related to Kosovo; and
- (f) social support services in communities, including services to children and care.

---

<sup>20</sup> Top Channel is a national commercial TV station located in Tirana.



**Figure 30: SP topic modelling Albania—Period 1**

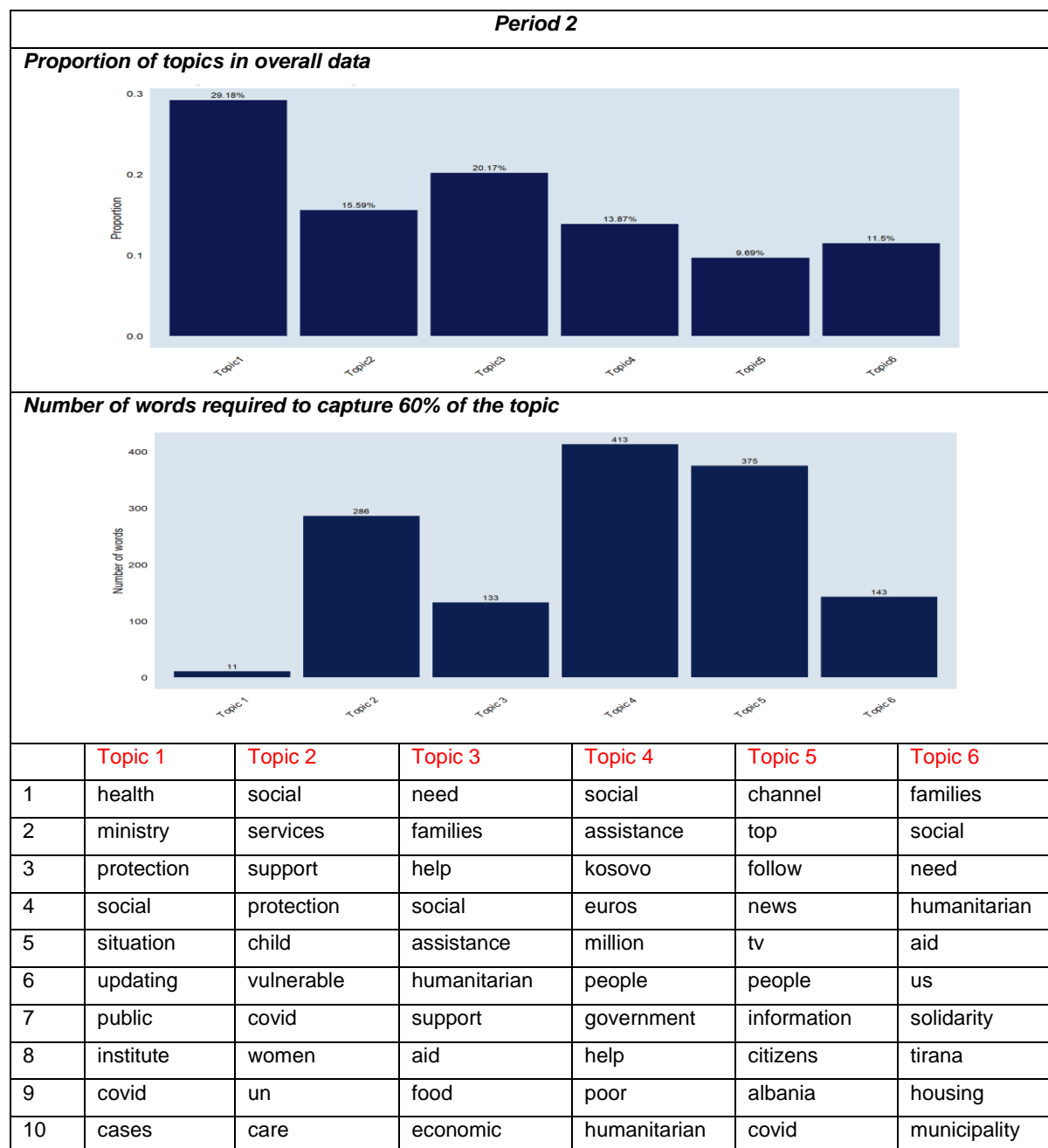


In Period 2, the topics were:

- (a) updates on the COVID-19 situation by the Ministry of SP and Health;
- (b) SP, services, and support for children and vulnerable households affected by COVID-19;
- (c) the provision of humanitarian and social assistance to families in need, including economic support and food aid;

- (d) the provision of social and humanitarian assistance for the poor in Kosovo;
- (e) Top Channel TV news and information on the COVID-19 situation in Albania; and
- (f) the provision of social and humanitarian aid by/to the municipality of Tirana related to housing.

**Figure 31: SP topic modelling Albania—Period 2**

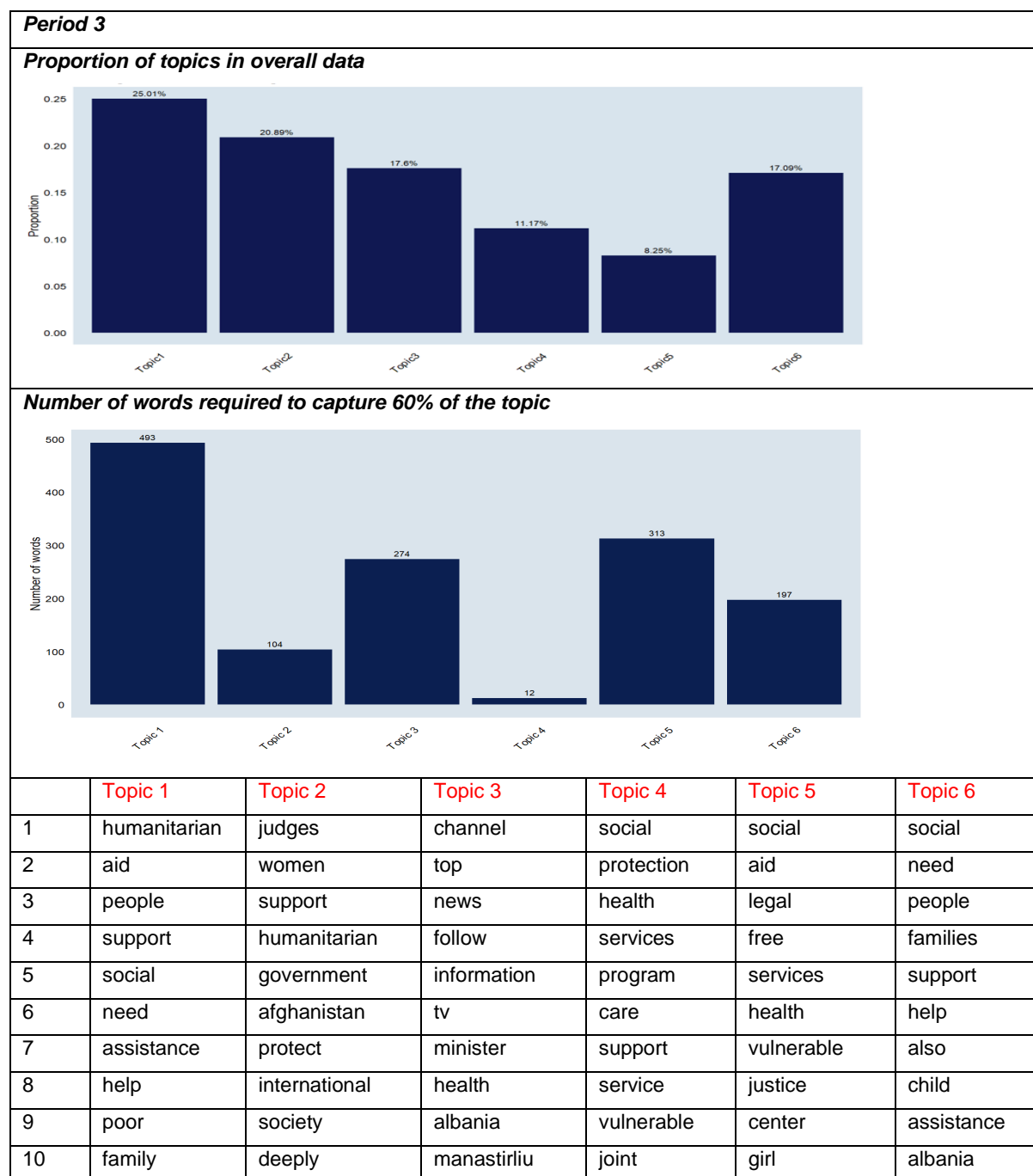


During the later COVID-19 period, the six topics identified by the LDA algorithm could be defined as follows:

- (a) the provision of humanitarian aid and social support to poor families;
- (b) issues related to justice with respect to women in Afghanistan;

- (c) Top Channel TV news regarding health-related topic discussed by the Minister of Health and parliamentarian Ogerta Manastirliu;
- (d) SP and health services and programmes for the vulnerable;
- (e) social and legal aid for the vulnerable, including girls; and
- (f) social support for people in need, including children.

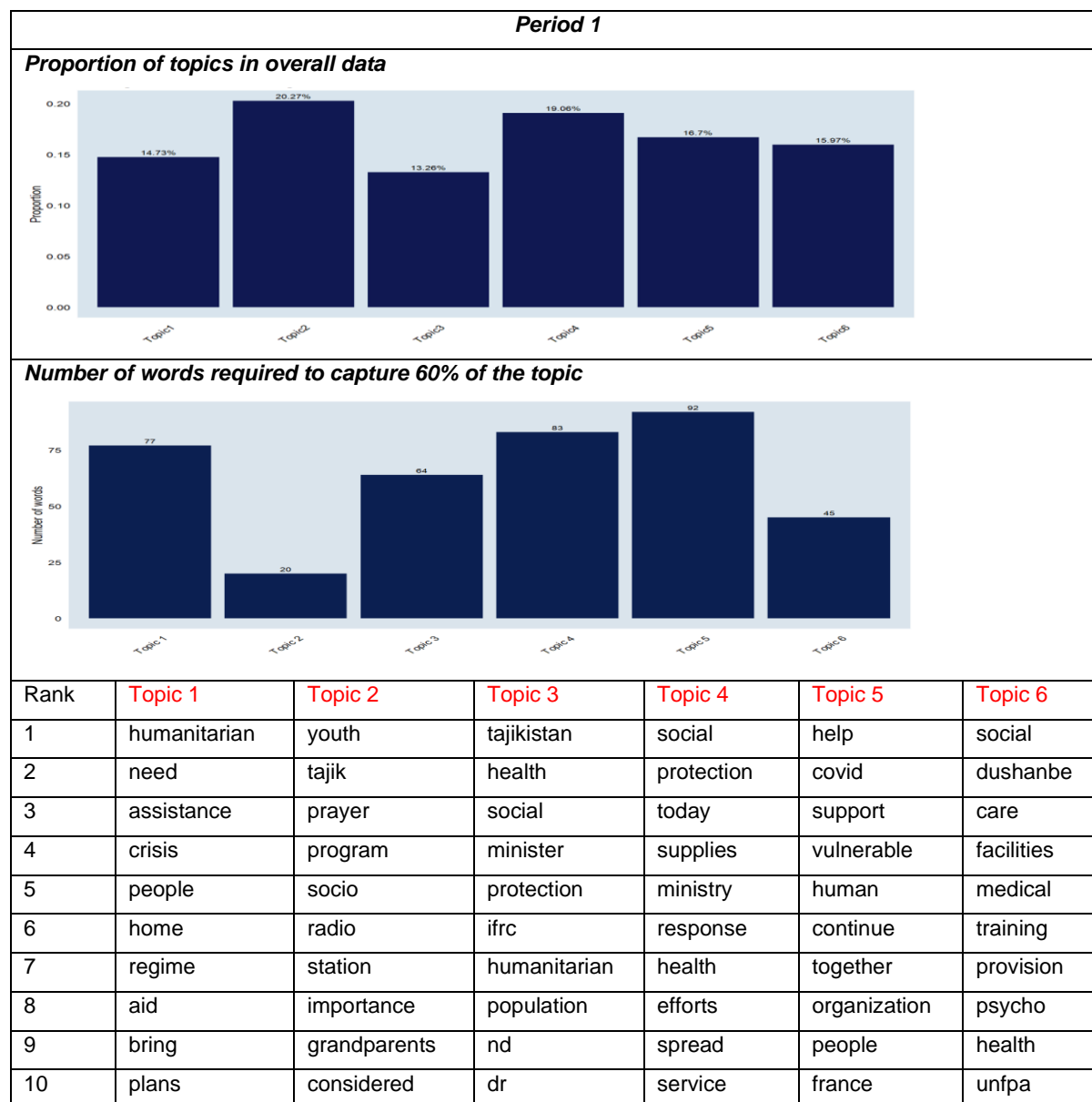
**Figure 32: SP topic modelling Albania—Period 3**



## Topic modelling detail for Tajikistan

We present similar detail for topic modelling on the Tajikistan text corpus below.

**Figure 33: SP topic modelling Tajikistan—Period 1**

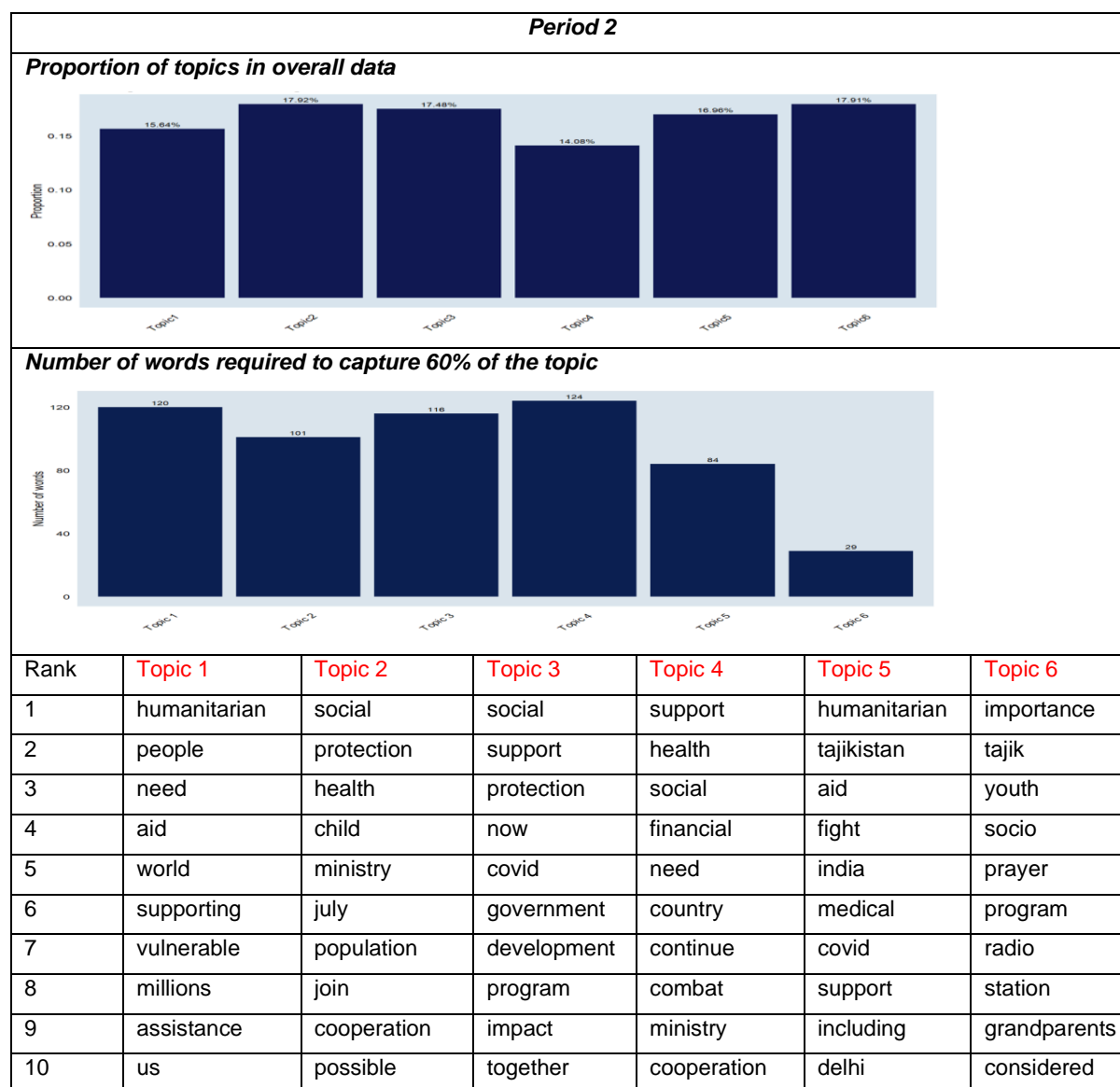


The topics above can be labelled as:

- (a) the provision of humanitarian need and assistance/aid in crisis;
- (b) traditional religious youth-centred radio programmes;
- (c) the humanitarian role of the Minister of Health and SP and IFRC;
- (d) the ministry's SP response and health supplies;
- (e) support to the vulnerable population affected by COVID-19 (mentions France); and

- (f) the provision of social care in Dushanbe, facilities, training, and attention to psychic health (mentions UNFPA).

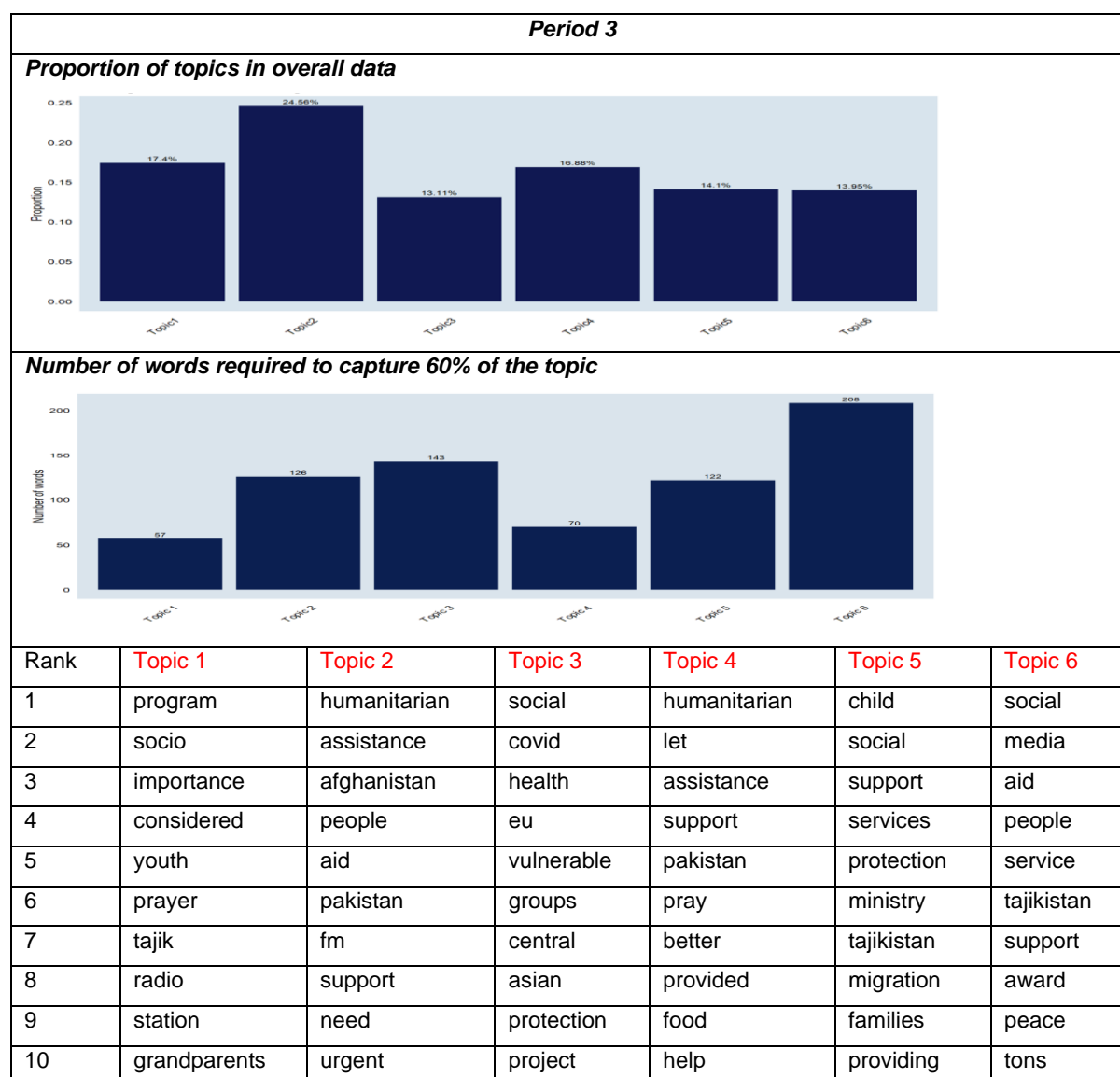
**Figure 34: SP topic modelling Tajikistan—Period 2**



The topics above could be labelled as follows:

- (a) the provision of humanitarian assistance/aid for vulnerable people around the world;
- (b) the Ministry of SP and Health: protection of children and cooperation;
- (c) the government's SP support and programmes linked to the impact of COVID-19;
- (d) financial support to address health and social needs and the ministry's cooperation;
- (e) the provision of humanitarian aid and medical support to fight against COVID-19 in Tajikistan and India; and
- (f) traditional religious youth-centred radio programmes.

**Figure 35: SP topic modelling Tajikistan—Period 3**



The topics above could be labelled as follows:

- (a) traditional religious youth-centred radio programmes;
- (b) the provision of humanitarian assistance in the face of urgent need (mentioning Afghanistan and Pakistan);
- (c) SP for **vulnerable** groups in connection with COVID-19 (mentioning European Union and Central Asian countries);
- (d) the provision of humanitarian assistance and support, including the provision of food (mentioning Pakistan);
- (e) SP and services provided by the ministry to children and migrants; and

- (f) social media coverage of aid to support peace efforts in Tajikistan as linked with events in Afghanistan.<sup>21</sup>

---

<sup>21</sup> To better label this topic, looking beyond the first 10 words was necessary because of the vague nature of the topic—it is the least concentrated of the six topics, judging by the number of words that cover 60% of the topic. Note that, during this period, the United States armed forces were withdrawing from Afghanistan and the Taliban was taking power, which threatened to destabilise the region.