



Inclusive  
Science and  
European  
Democracies



Grant Agreement No 960366

Project Acronym: *ISEED*

Project Title: *Inclusive Science and European Democracies*

Project website: [www.iseedeurope.eu](http://www.iseedeurope.eu)

Framework Programme H2020

Call H2020-SC6-GOVERNANCE-d under the H2020 programme

## **On Argumentation and Polarisation**

*Based on Deliverable D5.1*

*by*

*Giorgia Minello (UNIVE), Giuseppe A. Veltri (UNITN), Carlo R. M. A.  
Santagiustina (UNIVE), Massimo Warglien (UNIVE)*

### **Acknowledgements and disclaimer**

The ISEED project has received funding from the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement No 960366.

*This document reflects only the author's views and the European Union is not liable for any use that may be made of the information contained therein.*

# On Argumentation and Polarisation

The report on Argumentation and Polarisation is part of WP5 - 'Investigation of the prospects for deliberation using digital technologies' and focuses on the testing of the argument extractor tool. The argument extractor is a set of tools and methods for extracting and analysing causal statements and other forms of argumentation from social media and, more generally, online textual data.

The report is part of task 5.1.1 of WP5, which is intended to provide a set of tools that will be made accessible as web services via open APIs to retrieve and store web data from social media and other web platforms. In this report, first, we introduce some notion of NLP analysis applied in the context of social science; second, we present case studies where we tested the argument extractor tool about the IF-THEN relationship in two datasets of tweets about COVID-19 covering the three years of the pandemic and about climate change from 2017 to 2019. In particular:

- The first dataset used to test the Argument Extractor is called "*catch covid*". As the name suggests, it concerns how people debate (online) about catching coronavirus. The dataset has been built by downloading tweets from Twitter via the Tweet Downloader (fullArchive API endpoint V2); this API allows to download of large batches of Tweets into CSV or JSON files. In particular, we have downloaded 503 JSON files; 251 containing tweets and the remaining user data. The 251 JSON files collect in total 117753 tweets.
- The second employed dataset is named "lockdown", regarding online debates about the lockdown period. Even in this case, the dataset has been constructed by downloading tweets from Twitter, via the Tweet Downloader (fullArchive API endpoint V2. In particular, we have downloaded 51 JSON files containing in total 13000 tweets.
- A third dataset related to the "climate change" topic has been utilized as well. This last dataset is quite different from the previous two presented above as it has been not directly downloaded by Twitter Stream API but via its repository web page<sup>1</sup>, [Littman, Justin, and Laura Wrubel. "Climate Change Tweets Ids." <https://doi.org/10.7910/DVN/5QCCUU>, checked on 3.16 (2019): 2020], it was possible to download only tweet ids. The reason we could not get full-content tweets is that Twitter has recently restricted the redistribution of

---

<sup>1</sup> <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/5QCCUU>

Twitter content to third parties, so what we can get online are simply datasets consisting of tweet IDs of relevant tweets, also called dehydrated tweets. To obtain the tweet content, we had to hydrate these tweet IDs. The time window goes from September 21, 2017 to May 17, 2019, with a gap in data collection between January 7, 2019, and April 17, 2019. The original dataset contains the tweet ids of 39622026 tweets. However, it mostly retweets (namely thousand of duplicates) and/or single few words tweets. We processed less than 50% of them. For our analysis, we considered 2018-2019 tweets, filtering every single sentence based on the “if-then” mechanism. Eventually, our final climate change dataset is being consisted of 15000 sentences.

We have applied the argument extractor and some additional standard NLP techniques to explore the evolution of the online debate about COVID-19, its origins, the health policies and the vaccines. Similarly, we have applied this new tool to a corpus of data about climate change, the recent debate about policy initiatives, social movements and related issues. Both the word cloud and, above all, the word embedding representations, and in turn and the information obtained, did help us to reconstruct and understand (in broad outline) arguments and debates people had on online social platforms. The application of this tool represents a first layer of analysis that researchers can further explore. The polarization is, first of all, semantic in the sense that different implications and interpretations produce the formation of several sub-debates. It is not a surprise, given the nature of platforms like Twitter, that favour this type of process. While the focus is often on the issue of polarization as irreconcilable positions about one given issue, equally important are the proliferation of arguments leading to the creation of micro-debates that are insulated from each other. COVID-19 was the newest of the two topics analysed, the one in which little consolidated knowledge, and therefore argument, existed, and we can see the evolution of different ‘lines of argumentation’. The analysis of the climate change debate is less fragmented because there has been a debate now for years. The range of subtopics to discuss and, therefore, the associated arguments are less. The other consideration is that the polarization of arguments might have a distinct development pattern, and it should not be considered equivalent to the polarization dynamics of a debate.

The analysis in this report is at the macro level of analysis, in other words, a bird’s eye view. The use of NLP techniques, including a tool like the argument extractor, represents a useful instrument in the researcher toolbox. Still, it is insufficient to study a complex phenomenon like polarization that has to do with arguments as much as emotional responses and identity (Bail, 2021). The latter aspects have been discussed in the project in other deliverables (D5.2) and will be explored in the ISEED project.