

EGU21-8637

<https://doi.org/10.5194/egusphere-egu21-8637>

EGU General Assembly 2021

© Author(s) 2022. This work is distributed under the Creative Commons Attribution 4.0 License.



Combining Supervised and Unsupervised Learning to Detect and Semantically Aggregate Crisis-Related Twitter Content

Jan Bongard and Jens Kersten

German Aerospace Center, Institute of Data Science, Jena, Germany (jan.bongard@dlr.de)

The Twitter Stream API offers the possibility to develop (near) real-time methods and applications to detect and monitor impacts of crisis events and their changes over time. As demonstrated by various related research, the content of individual tweets or even entire thematic trends can be utilized to support disaster management, fill information gaps and augment results of satellite-based workflows as well as to extend and improve disaster management databases. Considering the sheer volume of incoming tweets, it is necessary to automatically identify the small number of crisis-relevant tweets and present them in a manageable way.

Current approaches for identifying crisis-related content focus on the use of supervised models that decide on the relevance of each tweet individually. Although supervised models can efficiently process the high number of incoming tweets, they have to be extensively pre-trained. Furthermore, the models do not capture the history of already processed messages. During a crisis, various and unique sub-events can occur that are likely to be not covered by the respective supervised model and its training data. Unsupervised learning offers both, to take into account tweets from the past, and a higher adaptive capability, which in turn allows a customization to the specific needs of different disasters. From a practical point of view, drawbacks of unsupervised methods are the higher computational costs and the potential need of user interaction for result interpretation.

In order to enhance the limited generalization capabilities of pre-trained models as well as to speed up and guide unsupervised learning, we propose a combination of both concepts. A successive clustering of incoming tweets allows to semantically aggregate the stream data, whereas pre-trained models allow to identify potentially crisis-relevant clusters. Besides the identification of potentially crisis-related content based on semantically aggregated clusters, this approach offers a sound foundation for visualizations, and further related tasks, like event detection as well as the extraction of detailed information about the temporal or spatial development of events.

Our work focuses on analyzing the entire freely available Twitter stream by combining an interval-based semantic clustering with an supervised machine learning model for identifying crisis-related messages. The stream is divided into intervals, e.g. of one hour, and each tweet is projected into a numerical vector by using state-of-the-art sentence embeddings. The embeddings are then grouped by a parametric Chinese Restaurant Process clustering. At the end of each interval, a pre-

trained feed-forward neural network decides whether a cluster contains crisis-related tweets. With a further developed concept of cluster chains and central centroids, crisis-related clusters of different intervals can be linked in a topic- and even subtopic-related manner.

Initial results show that the hybrid approach can significantly improve the results of pre-trained supervised methods. This is especially true for categories in which the supervised model could not be sufficiently pre-trained due to missing labels. In addition, the semantic clustering of tweets offers a flexible and customizable procedure, resulting in a practical summary of topic-specific stream content.