

Integrating Interactive and Computational Approaches for News Understanding

Liang Wu, Fred Morstatter, Daniel Baird, Justin Sampson, Kai Shu, Huan Liu

Arizona State University
Tempe, AZ, USA

1 Introduction

Event corpora such as GDELT¹ provide detailed summaries of events as they happen day-to-day. These datasets are unprecedented in their ability to extract machine-readable information from news data. However, this immense amount of data is nearly impossible to understand at the corpus level, requiring immense amounts of time in order to understand the meaning of, and relationships in between events [2]. The immense and rich data available in GDELT enables social scientists to ask questions previously unavailable. For example, the prominence of different actors, and the geographic distribution of different stories can be answered at a greater scale. While this analysis is suddenly possible, it is difficult to process this large amount of data. Reading the individual news stories, or even the meta data of them, is still impossible. Thus, we will integrate the power of topic modeling with the giant news source of GDELT to facilitate fast understanding of news stories.

Topic modeling [3] presents an unprecedented opportunity to take a large corpus of data and reduce it to a few core “topics”. These topics are often presented as lists of words which are then presented as a summary of the data. Topic modeling algorithms enable researchers and practitioners to quickly understand a huge corpus just by reading the few topics that are extracted from the corpus. The key issue in generating topics of interest is to identify the relevant data and suitable parameters. However, since domain experts, such as social scientists and policy makers, are not necessarily experts of computational tools, a new challenge arises for us as researchers to develop better ways to quickly understand event data by combining this rich data source with established methods like topic modeling.

In this demonstration we propose to integrate user interaction with computational tools, instead of generating results and then visualizing them. In particular, we will build a software tool which enables researchers to select/remove/modify topics that are being generated. The tool can help researchers and practitioners to quickly and accurately understand the immense data available in the GDELT dataset and find topics of interest effectively. Moreover, we will allow the users to specify both time intervals, as well as geographic and regional queries, and display the topics from the results of those queries.

¹ <http://gdeltproject.org/>

