Humanitarian Technology: Science, Systems and Global Impact 2015, HumTech2015

# Tackling the challenges of situational awareness extraction in Twitter with an adaptive approach

Haji Mohammad Saleem*, Faiyaz Al Zamal, Derek Ruths

*School of Computer Science, McGill University, Montreal, QC H3A 0G4, Canada*

## Abstract

Twitter is widely perceived as a potential source of valuable information for responders to mass emergencies. Despite interest in the development of extraction systems for such information, little effort has been put towards systemic methods for obtaining all posts pertaining to a disaster from the live Twitter stream. Researchers rely on keyword-based filters to extract information in spite of evidence that such markers are absent in many informational tweets, and also neglect the topic and traffic dynamics of the relevant tweets as crises progress. Previous work has shown that such practices can often lead to the loss of critical information in the context of a disaster. We introduce an adaptive filter, tailored to the idiosyncrasies of the real-time Twitter feed, intended to extract disaster-related content. Furthermore, we introduce a novel data model based on a three-label classification scheme to describe the composition of the data-stream. We use this model to simulate Twitter streams, modelling various post-disaster scenarios, for the purpose of filter performance evaluation. The filter is able to remove over 85% of the non-crisis content, and achieves a three-fold reduction in the loss of relevant contents compared to the existing approaches. In combination, the method and the model are useful tools for extracting situational awareness and highlight important directions for future work in this area.

## 1. Introduction

During and in the immediate aftermath of a major calamity, responders need to quickly assess conditions in the affected region. Such *situational awareness* (SA) enables the effective delivery of services and resources (e.g., humanitarian aid, and search and rescue operations) to the appropriate regions and populations [1]. However, disaster conditions impede assessment of affected locations and populations, making situational information difficult to obtain through traditional mean.

Social media has yielded a promising alternative source of such situational information: individuals in the affected regions, in many cases, readily post information about their conditions to platforms such as Twitter and Facebook. Such information may not be readily available and can contribute to SA. In fact, the analyses of recent disasters has shown that users of social media (notably Twitter) posted information would be valuable for first responders. [2,3]. As a recent example, immediately after the shooting incident in Moncton, New Brunswick, a Twitter user posted a

---

 * Corresponding author. haji.saleem@mail.mcgill.ca.

photo of the shooter, with a clear view of his armaments [1]. Such information, were it extracted rapidly and properly, could have been crucial in evaluating the threat posed by the shooter.

While Twitter data extraction is an active area of research (e.g., [4,5]), no approaches, to our knowledge, are tailored to the idiosyncrasies of disaster-related tweets in an active Twitter stream [7]. Recent work has established that disaster-related content in Twitter has several consistent characteristics:

*Small volume relative to the entire Twitter stream* - Since Twitter is used globally, an individual event generates only a small percentage of the overall traffic. This is true for disaster and non-disaster tweets alike [8].

*Absence of keyword labels* - A careful study of several disasters revealed that many early, information-laden tweets do not carry hashtags and other keywords which have been typically used for disaster-related tweet identification [9].

*Abundance of disaster tweets changes over time* - Though public response varies with disaster, the overall volume of resulting tweets can generally be split in to three phases: the *rise* demonstrating in-
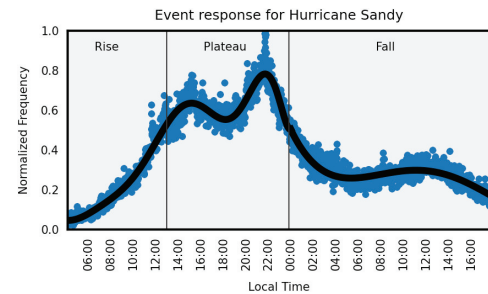


Fig. 1: *The three phases of event response on Twitter.* - The event response to Hurricane sandy on Twitter on Oct 29 and 30, 2012. The General shape can be divided into three sections, the initial rise, the middle plateau and the final fall, which is consistent with how the general population reacts to an event.

creasing frequency, the *plateau* demonstrating high frequency, and the *fall* demonstrating diminishing frequency of disaster related tweets, (see Figure 1) [8].

*Dramatic topical shifts* - Since an event is dynamic, the nature of public engagement changes as the event unfolds. For eg., the Moncton incident consisted of the initial shooting followed by an extended manhunt and finally the suspect's capture — all of which was reflected by discussion in Twitter. Moreover, many disasters, by their nature and coverage in the media, will engage larger Twitter populations well-outside the affected geographical area and generate additional content with diverse topics including emotional reactions, prayers for victims, and offers of support.

The need to account for such issues motivates the two major contributions of our study: (1) an adaptive filter that can accommodate for the overwhelming amount of unrelated content in the stream along with the variable volume and dynamic topical shifts of relevant content over time, and (2) a formal data model describing the composition of an active Twitter stream in terms of the relative frequencies of different types of disaster-related content.

### Adaptive filter

Twitter generates almost 500 million tweets everyday. Therefore, any extraction system that operates on the entire data-stream has to face excessive noise in the form of non-relevant content. If not addressed carefully, this creates a learning bias towards the more dominant label, degrading the overall performance [10,11]. We propose a noise filter as an essential part of any information extraction machinery to eliminate data imbalance by effectively reducing the non-relevant content in the data-stream.

To account for the dynamic nature of the event stream, both in volume and content, we develop the filter with an adaptive framework that periodically updates its model based on recently labelled data. Actual noise filtering is carried out by a supervised learning algorithm embedded in the adaptive framework. The choice of algorithm depends on the following factors: (1) we require a low complexity algorithm to reduce computational overheads due to the real-time nature of data processing pipeline to which the filter belongs; and (2) we need an algorithm that performs reasonably well with sparse datasets due to the short document size of a single tweet. We therefore embed a Naive Bayes classifier in our filter frame work, appropriately tailored to deal with the adaptive requirements and adjusted to handle the data imbalance present in the Twitter stream.

We require pre-labelled data to test our supervised learning setup and properly measure its performance. However, labelling even a small timeframe of an entire Twitter stream is not logistically possible. To overcome this problem, we generate synthetic datasets that represent simulated Twitter streams that stem from the data model we soon explain. The adaptive Naive Bayes filter successfully removes the majority of the noise in these simulated streams, while losing only a small fraction of relevant data. Our approach outperforms keyword-based filters, non-adaptive filters or

---

[1] https://Twitter.com/PatHemsworth/status/474356870686461953