# Using Twitter to Detect and Investigate Disease Outbreaks

**David Marchette\* and Elizabeth Hohman**

Naval Surface Warfare Center, Dahlgren, VA, USA

## Objective

In this work we investigate the extent to which social media, in particular Twitter, can be used to detect an outbreak of a disease or illness. We term these outbreaks "events", and we will describe methodologies for detecting events.

## Introduction

Social media is of considerable interest as a sensor into the thoughts, interests and health of a population. We consider three types of health events that an analyst may wish to be made aware of:

- Given a known disease, such as MERS, SARS, Measles, etc., an event corresponds to individuals contracting the disease.

- Given a set of symptoms (fever, stomach pain, etc.), an event is an unusual number of individuals1 complaining of the symptoms.

- Most generally: an event is an unusually large group of individuals who can be identified as being effected by some personal illness.

Note that to detect an "unusual number" of something, we need to count the indicators of the event, and we need to compare the current count with past counts. Further, we are generally interested in geographically constrained events, and so for this work we will focus on county-based counts. We will count the number of items (tweets or individuals) expressing the event indicator (a disease name, symptom, or classified as "personal health related" as indicated by our classifier). Our approach to detecting health related events is: filter -> classify -> detect. We first filter out tweets that contain no "health related" terms, then apply a classifier to each tweet. This classifier is designed to flag a tweet as being about "personal health" or not. We then aggregate the positive instances per day at the county level and detect as an event any county/day pair with an unusually high count (as compared to the recent past).

## Methods

We collect all tweets with latitude and longitude within the continental United States. These are then filtered with a set of phrases which are designed to retain tweets that might be related to some disease, sickness, or other health event. These phrases are symptoms (fever, headache, feel sick, etc.), disease or pathogen names (flu, a cold, salmonella) and remedies (cold meds, ibuprofen, Nyquil). A random forest, trained on 13K hand tagged tweets, is then applied to the matching tweets to classify them as about "personal health" (class 1) or not (class 0). Finally, we count the number of class 1 users within a county per day, and compare this to the past using a sliding window z-score approach. Any county/day pair with an unusual count (over 3 standard deviations) is flagged as an event and is then assessed by the analyst as to its nature and importance.

## Results

We illustrate the methods through detections of the Boston Marathon bombing and several events related to New Year's festivities. We also show how discussion of the Ebola outbreak in West Africa changes over time.

## Conclusions

The approach of filter -> classify -> detect, wherein we first consider only those tweets matching certain "health related" terms, then classify the tweets as being about "personal health" or not, and finally detect anomalies based on localized counts, is a very powerful method for processing these data.

## Keywords

Biosurveillance; Social media; Twitter; disease outbreak; text analysis

**\*David Marchette**
E-mail: david.marchette@navy.mil