# Data Science for Social Media in Emergency Management

Interim Report for NM Supercomputing Challenge 2017

**Team Number:** PESSF118
**School Name:** Piñon Elementary School, Santa Fe
**Area of Science**: Data Science

## Problem Definition:

People communicate in multiple ways during disasters: face-to-face, by watching TV, and by using social media.  Social media is rapid and two-way, but it lacks context and is always changing[1]. **Social Media in Emergency Management (SMEM)** aims to improve communication among the public, first responders, and emergency managers.  SMEM is new and has only recently been studied.[2-6] Our team's goal is to **use data science to improve SMEM.**

SMEM work is done by a Virtual Operations Support Team (VOST) which reports to an Emergency Manager.  We had a VOST staffer, Marlita Reddy-Hjelmfelt (who leads the Pacific Northwest VOST), visit us and explain VOST operations and needs.  At present, VOST work is usually done by volunteers with little or no budget, with no special data access, and little software support.  Ms. Reddy-Hjelmfelt described the main issues for VOST workers as:

1.  **Finding good search terms to retrieve relevant posts.**  Many posts do not include standard hashtags.  Few posts are geolocated.
2.  **Doing complex searches on streams**.  VOST workers rely on a combination of mental filtering and personal requests to friendly academics".
3.  **Distinguishing false from true reports.**  An emergency attracts some people who falsify reports, which have to be distinguished from real reports using human judgement.

Ms. Reddy-Hjelmfelt did not have a high opinion of machine-learning classification of social media streams.

**Our hypothesis is that we can use artificial intelligence on deep Twitter streams containing photographs to tell pet owners from non-owners.**  We chose this problem to stand in for the problem of distinguishing fake and real reports because most people do not have privacy concerns about saying whether or not they own pets.   This question also has practical importance in SMEM. We hope that the tool we build can be used for routine future VOST use.

## Problem Solution:

We propose to build a tool that classifies social media streams based on both boolean searches and machine learning using Twitter data.  We propose two new features for our classification engine: (1) use of text tags from image analysis as part of the data, and (2) searching deep into Twitter feeds rather than by doing classification on single tweets.  We will implement this approach by hacking on code from the example "[AWS re:Invent 2015: real-world smart applications with Amazon Machine Learning"](#).  This pipeline was designed for classification of customer service complaints using simple serverless python scripts.  Serverless python can be scaled to data supercomputing scale using the Amazon Lambda service.  The three models that we propose to investigate are:

1. Boolean-search classification on deep streams of Tweets with image classification.
2. Machine learning on deep streams Tweets without image classification.
3. Machine learning on deep streams of Tweets with image classification.
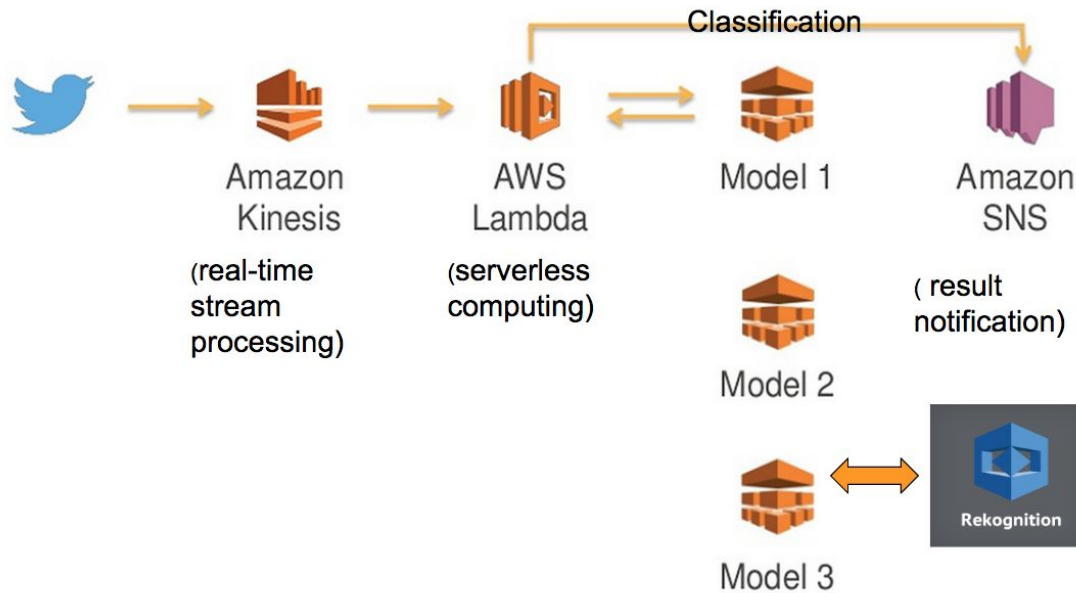


**Figure 1:** Proposed pipeline.

We propose to obtain the classification (pet-owner/non-pet-owner) for **1000 Twitter accounts** in using an entertainment service we built that matches user photos with celebrities, called "funyun". We built funyun just to test our code, but it proved to be quite entertaining and is getting multiple hits per day.  The entertainment service, which we call "funyun", identifies the closest matching to celebrity to user-supplied photos.  We will promote use of the service among the Piñon community and Santa Fe Public Schools in general through articles in take-home newsletters and probably a press release. We plan to divide the 1000 Twitter accounts into a group of 900 accounts to be used for training and 100 accounts to be used for testing accuracy, which will give a **10% noise level in characterizing classification**.  We will select the 100 test accounts randomly, but with an even distribution between the two classes to simply analysis (about 68% of households have pets).

Progress to Date:

We originally proposed to predict the location of pets in disasters using geolocated social media posts, but we were forced to pivot from this problem due to lack of data.  We developed a formula for estimating the sensitivity of social media for emergency management problems:

$$\text{Sens} \cong (\text{MAU/Pop}) \times \text{FracRel} \times \text{FracPub} \times \text{FracGeo,} \qquad \text{(Eq. 1)}$$

where *Sens* is the sensitivity, *MAU* is the number of users, *Pop* is the size of the relevant population (for example, the US population of people is 326 million),  *FracRel* is the fraction of posts that are relevant to the problem, *FracPub* is the fraction of user accounts that are open to public searching, and *FracGeo* is the fraction of social media streams that can be geolocated.  The following table shows approximate values for Equation 1 for the three social media sites that are most useful:

| Platform | US MAU, Millions | US MAU/ Pop | Frac Pub | Main Uses | Frac Rel, % | Frac Geo |
|---|---|---|---|---|---|---|
| Facebook | 214 | 65% | 20% | Connecting with friends | medium | >1% |
| Instagram | 150 | 46% | 32% | Pics of kids and pets | low | >1% |
| Twitter | 68 | 21% | 88% | Personal news | high | >1% |

**Table 1:** Sensitivity of some social media platforms for classification.

The problem is with the final column, *FracGeo;* very few posts on any media are geotagged. GPS locations on photos produced by smartphones (in EXIF data) are stripped by all relevant social media platforms. The other columns of Table 1 also shows why Twitter is the best source due to its public and news-like nature.

We have used good software practices in implementing *funyun.* The code lives in a source code repository at https://github.com/EagleBytes2017/funyun, and we have used appropriate software engineering tools (pypi, Travis, codecov, pyup, and codacy).

Expected Results:

We will test our classification tool on the problem of classifying Twitter accounts into pet owners and non-pet owners. We will characterize performance of our classifier using Receiver Operating Characteristic (ROC) curves for the three automated models plus human classification by the team.

Team:

**Members:** (Aliases used because of privacy restrictions) DarkDJ, Assistanceskaterdog, soccerchamp, spursqueen volleyballqueen, Warmachine

**Sponsoring Teacher:** Delara Sharma.

**Mentor:** Joel Berendzen (GenerisBio, LLC)

References:

1. Petronzio, Matt (2017) "Facebook opens its Community Help API to disaster organizations". *Mashable*, November 29, 2017.
   http://mashable.com/2017/11/29/facebook-community-help-api-fundraising
2. Zagorecki, Adam & Johnson, David & Ristvej, Jozef. (2013) "Data mining and machine learning in the context of disaster and crisis management". *Int. J. Emerg. Mgmt.* **9**.351-365. DOI:10.1504/IJEM.2013.059879.
3. Ilyas, Andrew. (2014) "MicroFilters: Harnessing twitter for disaster management". Published in *Global Humanitarian Technology Conference*, 2014 IEEE. DOI:10.1109/GHTC.2014.6970316
4. Bruns, Axel & Burgess, Jean (2014), "Crisis Communication in Natural Disasters: The Queensland Flood and Christchurch Earthquakes". In Weller, Katrin, Bruns, Axel, Burgess, Jean, Mahrt, Merja, & Puschmann, Cornelius (Eds.) *Twitter and Society* Peter Lang, New York, pp. 373-384. http://eprints.qut.edu.au/66329
5. Smith, Brian (2010), "Socially distributing public relations: Twitter, Haiti, and interactivity in social media". *Public Relations Rev.* **36**:329-335. DOI: 10.1016/j.pubrev.2010.08.005