# AI & SOCIAL MEDIA

## based incident notification system

By Jiahua Gong
December 3, 2019

# 1 Introduction

Social media is not only a way to help people connect, but can also produce valuable data that makes people's life better.

The nature of social media that users generate the content, in a real-time manner, makes it a perfect source for incident detection, such as earthquake. However, there is challenge on analyzing and categorizing raw data.

This project aims to utilize AI to extract incident info from geo-tagged tweets and notify people nearby that may be affected to reduce losses.



*Figure 1. Real time incident map*

# 2 Method

Unsupervised learning and supervised learning are both applied to solve this problem.

First step is to train a classifier model using training dataset, i.e. tweets. The output is a list of weights for each word.

Second step is to run K-Mean algorithm to cluster the testing dataset by latitude and longitude.

Third step is to iterate each cluster and run classifier against each tweet to compute the relevance to an incident and overall percentage in that cluster to determine whether there is an incident.



*Figure 2. Tweets during Hurrican Sandy*

# 3 Results

**1. Train classifier model**

Input 104271 tweets that were geo-tagged and labeled.

For each tweet, extract features and compute the projected value and hinge loss.

Repeat 20 iterations with eta = 0.01, and then get a trained model which contains a list of weights for 102652 words.

**2. K-mean cluster**

Aggregate testing dataset which has 798 tweets based on longitude and latitude.

Start clustering from K=1, and then increase K by 1 for each re-cluster, until new centroid is close to an existing one.

Here when K = 3, centroids locate in New York City, Massachusetts, and Washington DC. And when K=4, new centroid locates again in New York City and clustering stops here.

**3. Analyze testing dataset of tweets**

Cluster 1(Elizabeth, NJ 07201):
a. Hurricane relevant: 535
b. Hurricane not relevant: 1
c. Hurricane possibility: 99.26%

Cluster 2(Garrett County, MD 21520):
a. Hurricane relevant: 106
b. Hurricane not relevant: 1
c. Hurricane possibility: 99.07%

Cluster 3:(Milford, MA 01757):
a. Hurricane relevant: 41
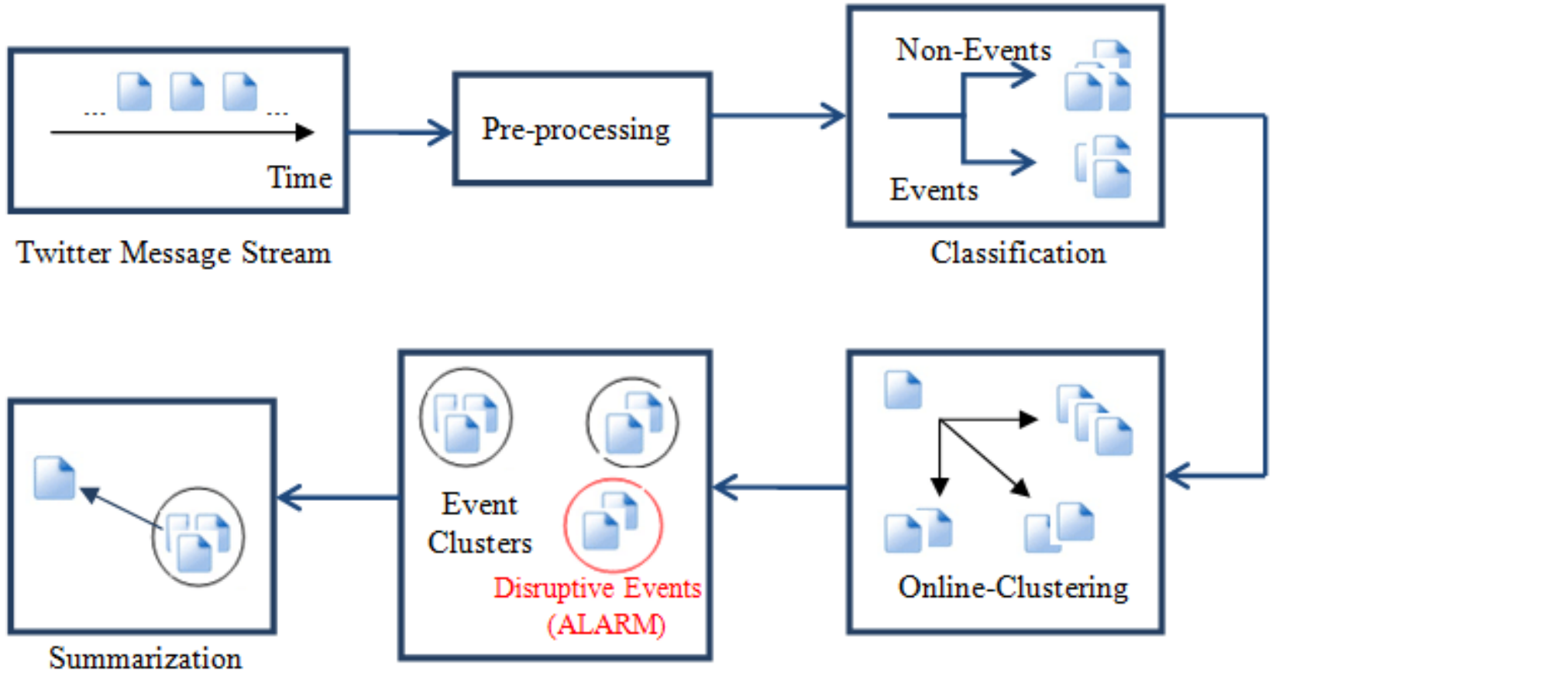b. Hurricane not relevant: 2
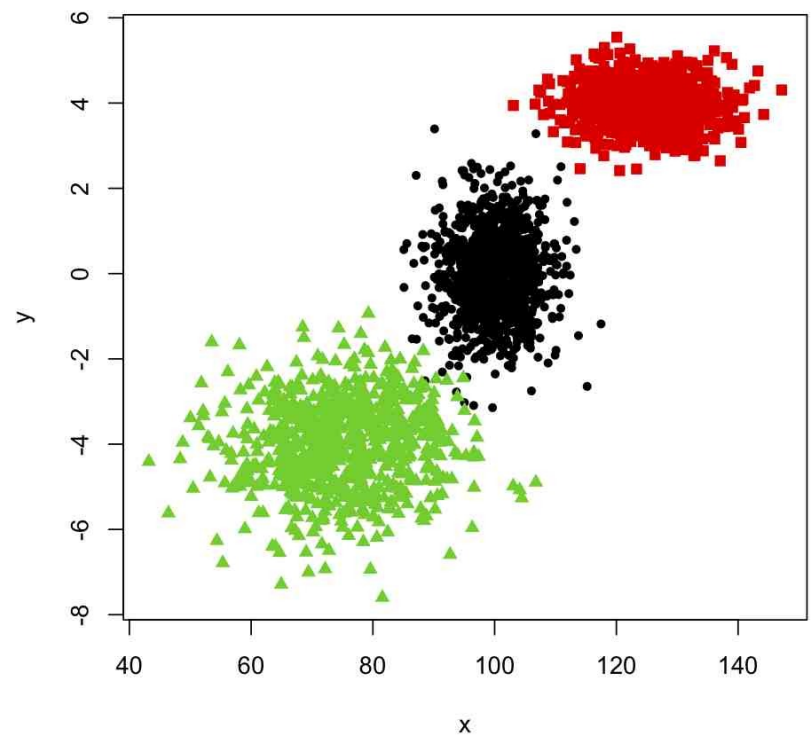c. Hurricane possibility: 95.35%



*Figure 3. Workflow*



*Figure 4. K-mean clustering*
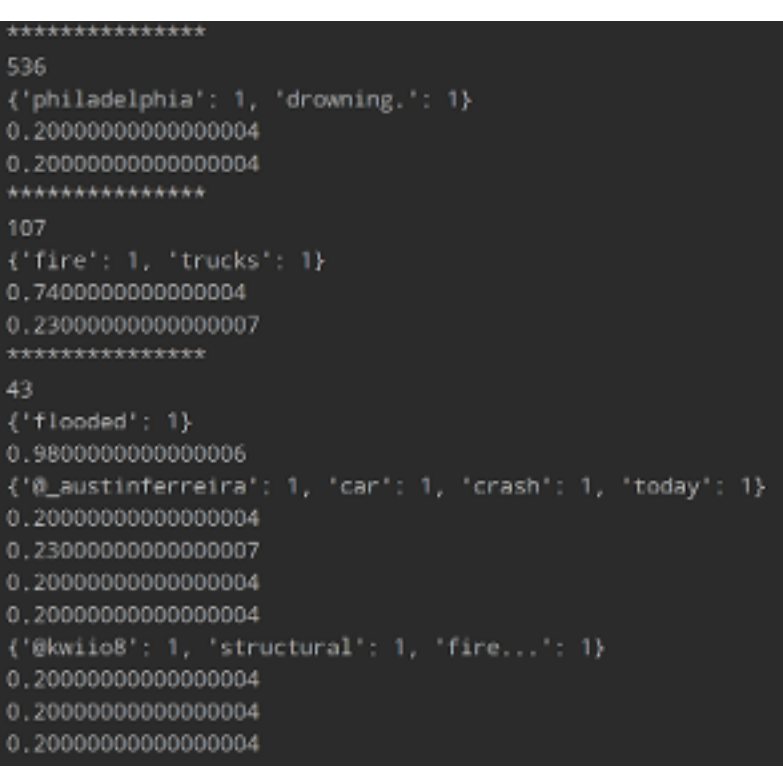


*Figure 5. Dataset*



*Figure 6. Result*

# 4 Conclusion

The result is very promising that not only the confidence of incident, e.g. hurricane, is very high, but also it actually shows the path of this hurricane.

The model is built upon statistic frequency and relevance which proved to be more accurate than keyword frequency or indirect indicators such as tweet length, sentiment, or punctuation. Also, with over 100k training tweets for this supervised learning process, the model demonstrates a very high effectiveness.

Unsupervised learning, i.e. clustering, shows a clear path of this hurricane. However, given the small size of testing dataset which are 798 tweets, future work needs to be done to dynamically prioritize the clusters to skip un-categorized or non-significant ones which will also improve the system performance in general.

Overall, the project turns out to be a successful experiment to utilize AI and social media data to detect real time incident aggregated by geolocation.
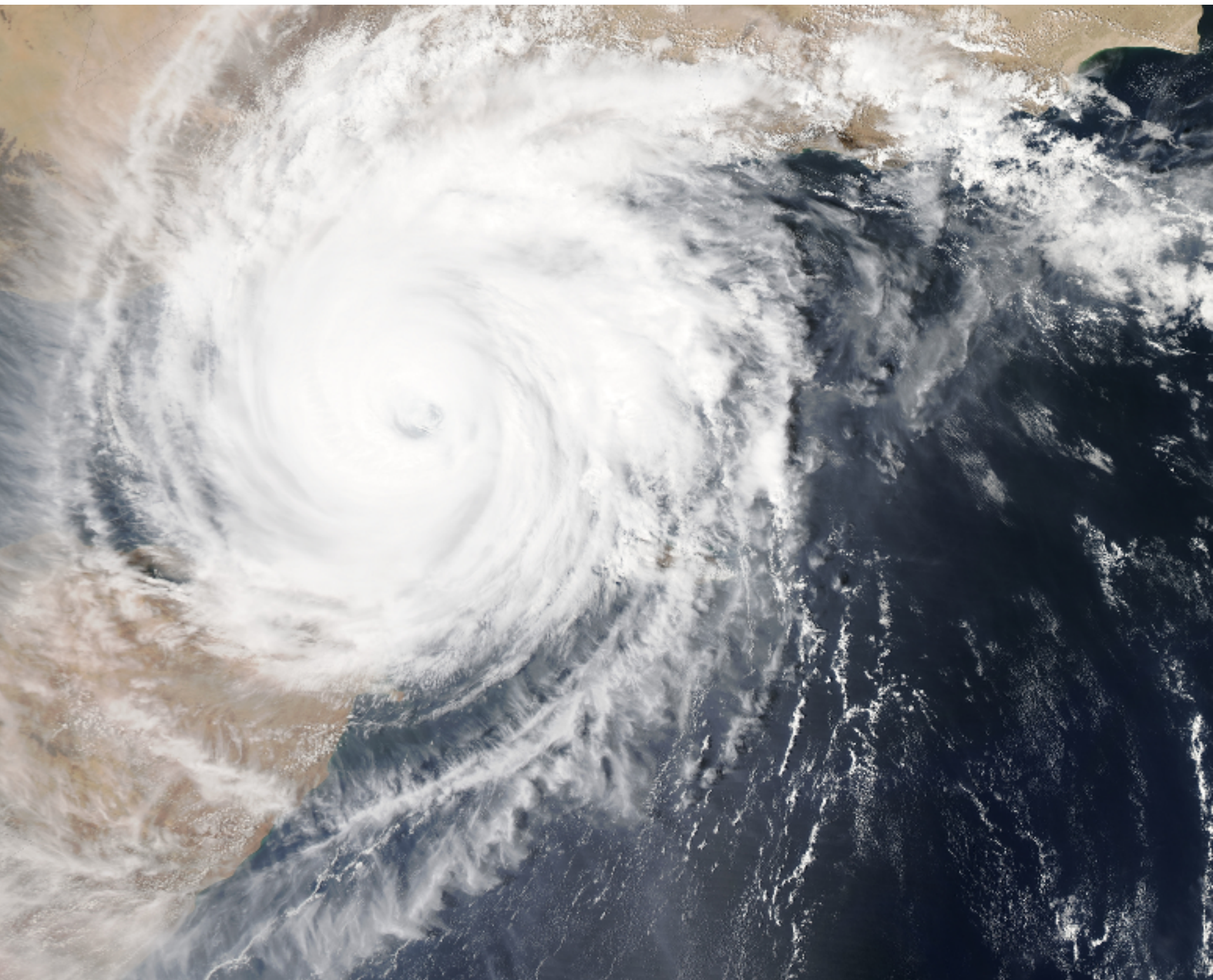


*Figure 7. Hurricane*



*Figure 8. Hurricane Alert*