



Leveraging Twitter to Inform Disaster Response Time

Jillian Berman,
Sree Gollakota, Zhan Yu





Problem Statement

In order to better improve emergency response times during disasters, this project aims to leverage social media to identify real time road closures or damaged roads, and other blocked routes that may affect travel time, traffic lights, safety, etc.

Case Study: Hurricane Michael



- Made landfall October 10th, 2018
- 74 deaths were attributed to the storm.
- Caused an estimated \$25.1 billion USD in damages.



Workflow/ Methodology



- 1.) Data Collection:
 - Scraping tweets from GetOldTweets3 library
- 2.) Data Analysis/Preprocessing:
 - NLP cleaning
- 3.) Classification between traffic incident tweets and normal tweets:
 - Use old tweets to train model for live classification
- 4.) Pull live tweets to feed into the classification model:
 - Scraping live tweets with Tweepy library and Twitter API
 - Run classification model on live tweets
- 5.) Map emergency tweet locations

Data Collection



- Used GetOldTweets3 Python library
- Date Range:
 - 10/01/2018 - 10/15/2018
- Twitter accounts scraped:
 - @BayCountyTMC , @fl511_panhandl & @WJHG_TV
- Keyword query:
 - 'HurricaneMichael' and '850strong' for example
- Training dataset contained 4,350 tweets.

Gathering Live Tweets



- Using the Tweepy Python library and a Twitter API we were able to scrape live tweets for the past 7 days.
- Scrapped the usernames and the same keywords we'd scrapped from our historical data, we were able to gather a dataset of 1014 tweets.

Data Analysis / Preprocessing



Cleaning tweets:

- Remove website links
- Remove non-letters
- Convert to lower case
- Remove basic stopwords

Options:

- Customize more stopwords
- Stemming/Lemmatizing

Modeling



- Corpus Vectorizers
 - CountVectorizer() or TfidfVectorizer()
 - Stop_words: English
 - Monogram, bigram, and trigram ranges
- Models
 - Random Forest, AdaBoost Classifier, Support Vector Machine for Classification
 - Choose best model by comparing training, and testing accuracy score
- Validation
 - Confusion matrix: test for optimized sensitivity

Modeling



- Best model #1: Count-vectorized AdaBoost Model
- Low bias and low variance in accuracy scores
- Outperforms baseline

<u>Best Parameters</u>		<u>Best Train Score</u>	<u>Best Test Score</u>	<u>Baseline</u>
<u>cvec_ngram</u>	(1,1)	0.99	0.973	0.841
<u>range</u>				
<u>stop_words</u>	None			

Modeling



- Best model #2: Count-vectorized Support Vector Machine for Classification
- Low bias but slightly larger variance than AdaBoost model

<u>Best Parameters</u>		<u>Best Train Score</u>	<u>Best Test Score</u>	<u>Baseline</u>
<u>cvec_ngram</u> <u>range</u>	(1,1)	0.995	0.951	0.841
<u>stop_words</u>	None			

Modeling



- Should be able to predict better for false negatives
- Compare sensitivity of AdaBoost model and Support Vector Machine

<u>Count</u> <u>Vectorized</u> <u>AdaBoost</u>	<u>Actual</u> <u>Positives</u>	<u>Actual</u> <u>Negatives</u>	<u>Sensitivity</u>	<u>Specificity</u>
<u>Predicted</u> <u>Positives</u>	True Positives 148	False Positives 4	0.86	.996
<u>Predicted</u> <u>Negatives</u>	False Negatives 25	True Negatives 909		

Modeling



- Should be able to predict better for false negatives
- Compare sensitivity of AdaBoost model and Support Vector Machine

<u>Count</u> <u>Vectorized</u> <u>SVM for</u> <u>Class</u>	<u>Actual</u> <u>Positives</u>	<u>Actual</u> <u>Negatives</u>	<u>Sensitivity</u>	<u>Specificity</u>
<u>Predicted</u> <u>Positives</u>	True Positives 121	False Positives 1	0.699	.999
<u>Predicted</u> <u>Negatives</u>	False Negatives 52	True Negatives 912		

Modeling



- Count-vectorized AdaBoost is more sensitive and will be used for mapping
 - Minimizes false positives
 - Better able to recognize emergencies and traffic incidents without ignoring true emergencies and traffic incidents
- Model limitations:
 - Model actually classifies between tweets from specific accounts and tweets without accounts specified
 - May not perform well with live tweets

Extracting locations with spaCy



- Using the spaCy Python packages
- Training our own spaCy model by providing many examples to meaningfully improve the system, since the default spaCy model performs poorly.

New: Object on roadway in Okaloosa TWN on I-10 west INS at MM 51 MM ,
right lane LN blocked. Last updated at 04:45:39PM. #f1511

- With trained spaCy model, we were able to extract locations from tweets.

Getting Latitude & Longitude



- HERE.com provides mapping and location data services.
- Using the herepy python package and HERE Geocoder API we were able to get latitudes and longitudes from most of addresses we extracted from historical tweets.

Mapping



- Using shapley, GeoPandas and folium we plotted the points we were able to generate from spaCy processing of the tweets.
- [Map Demonstration](#)

Conclusions and Limitations



- Gathering tweets
 - Model will perform better if classification assumptions of gathered training data are more robust
- spaCy processing
 - spaCy performs better when we train our model with more examples
 - Training the model is time intensive though
- Coordinate extraction
 - From our training data set we were able plot 20% of the tweets.

Next Steps



- From here if you sign up for a Google Maps API you can integrate live traffic conditions, as well as point to point instructions.
- Improving classification assumptions to be more specific to the tweet context as opposed to generalizing 511 account tweets as positive class and others as negative class

A polar map of Antarctica, showing the continent and surrounding islands. The map features a grid of latitude and longitude lines. The text 'Thank you.' is prominently displayed in the center. Various geographical features and historical sites are labeled, including King Edward VII Land, Victoria, and several islands like Alexander I Ld and Kaiser Wilhelm II Ld. The map also shows the names of explorers and dates, such as Cook 1774 and Scott's Earliest 1902. The word 'F' is visible in the upper and lower central regions of the map.

Thank you.