# LEVERAGING SOCIAL MEDIA TO MAP DISASTERS

CAMERON BRONSTEIN | DAN KIM | PATRICK O'CONNELL

# Agenda

1. The problem
2. Data collection + methodology
3. The modeling process
4. Conclusion + recommended next steps
5. Q&A

# Problem Statement

**How can we leverage social media messages to map natural disasters?**

Goal: Map specific areas that need assistance

1. Identify people that need help
2. Pinpoint their specific location

# Data & Methodology

**Third-Party Dataset**

- Disaster related messages
  - Included News Reports, Social Media, and transcribed phone calls with world health organizations.
- Caveat: Did not include geolocation data

**Twitter API**

- Search by location and time
- Allows us to classify tweets that have mappable locations

# Model 1

**Goal**: Identify between people in urgent need of help after a natural disaster and media that are referencing the event.

**Data**: Figure Eight - Multilingual Disaster Response Messages

- Included text from news media, social media, and direct communication between world health organizations and survivors on the ground.
- Direct messages:
  - Mostly included calls for help and resources
  - Used as a proxy for new "Urgent" class label
- News and Social Media
  - Used as a proxy for "Non-Urgent"

# Model 1

**Binary Classification Problem**

- "Urgent" vs. "Non-Urgent" Messages
- Model: Logistic Regression
  - High Model Accuracy, Interpretable Process

**Natural Language Processing**

- Count Vectorizer
  - one and two word combinations (n-grams); 50,000 features.
  - **Common "urgent" terms:** 'information', 'note', 'am', 'message', 'digicel', 'like to', 'what', 'find', 'give', 'street', 'please', 'send', 'you'
  - **Common "non-urgent" terms:** 'santiago', 'sandy', 'http', 'concepcion', 'haiti earthquake', 'rt', '000', 'frankenstorm', 'hurricanesandy', 'also', 'hurricane', 'donation', 'storm'

# Model 1

**Model Performance:**

| Accuracy: 0.974 | Sensitivity: 0.968 |
|---|---|
| **Misclassification Rate:** 0.026 | **Specificity:** 0.977 |

**False Negatives:**

*"When I got home, I saw the furniture moving backward. Then, the wall in the corner fell on top of it with all the displays and the television. The problem we have we need water and hunger is killing us. .. "*

*"My pregnant wife is in the hospital St.Therese in Miragoane. On the 12th she was hit by debris. She was told that the foetus has moved. I have no money. "*

# Model 1

**Further Model Validation**

- Test Data from Kaggle "Disasters on Social Media" dataset

**Predicted Urgent** - contains "please", "safe", "you"

- *"Thinking of our friends in Texas as they prepare for Hurricane Harvey.
  Please keep safe - we are all praying for you all to be OK!"*

**Predicted Non-Urgent -** contains news information and hyperlink

- *"Texas is bracing for extreme weather as Hurricane Harvey strengthens
  http://nyti.ms/2v9BCEK"*

# Model 1 Limitations

- Without geolocation, our predictions are not actionable.

- Searching broadly across platforms or search terms is inefficient

- "Proof of Concept"
  - Not necessarily generalizable across regions or events, even when trained on a broad corpus of data.

- Trained on considerable noise since labels were proxy.

- Future models will rely on more reliably collected and labeled data.

# Model 2 Overview

**Plan:**
- Proof of concept for modelling methodology
- Does not require classifying each individual tweet
- Geographical area (e.g. zip code) & time period = individual observations

**Goal**:
- Identify people in need of help due to, for example, flooding
- Get historical data including a natural disaster to train the model
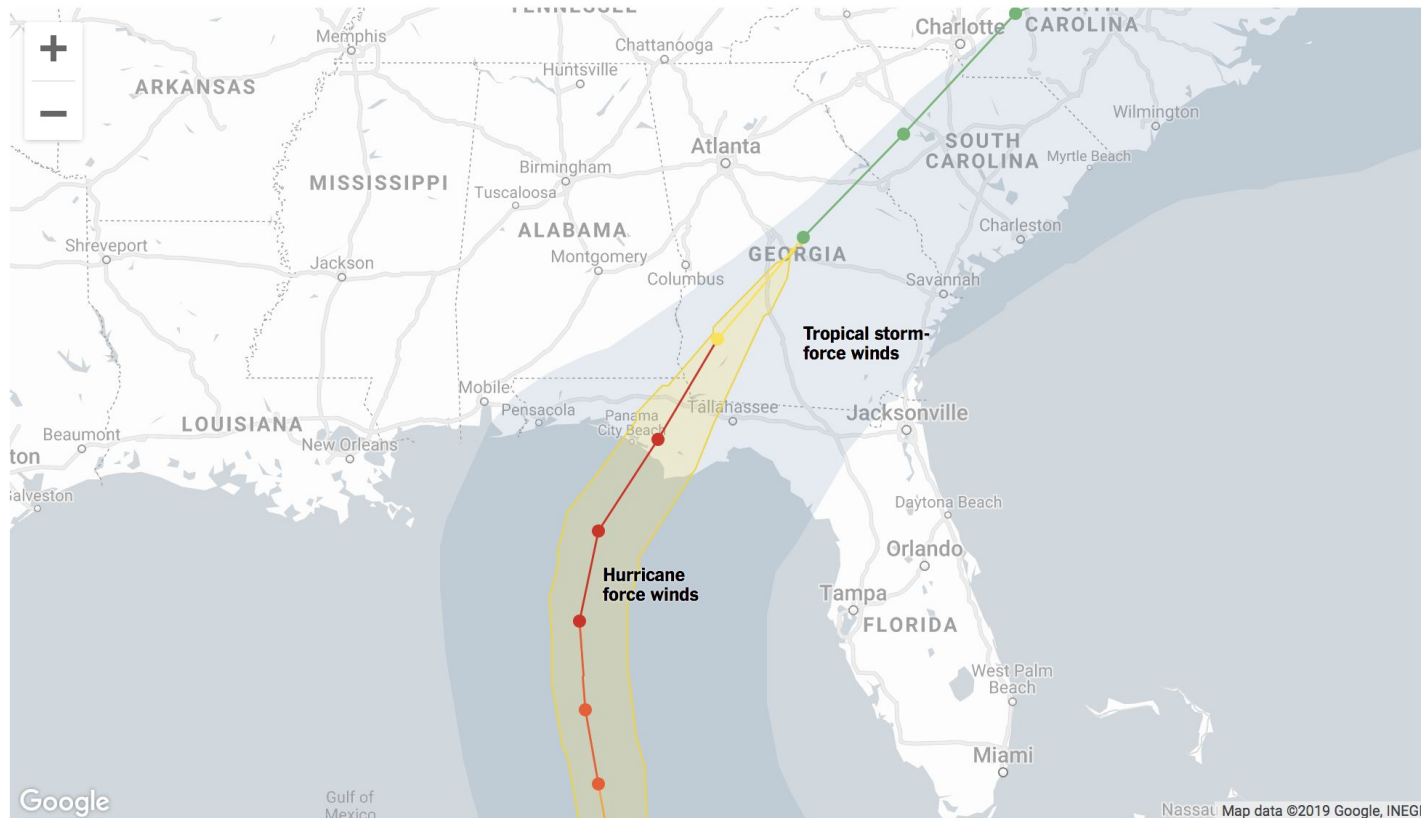- Target/label for each observation would be presence of natural disaster

*The model would then extract the important words from the tweets for that time and place, that indicated whether there was flooding (etc) or not.*
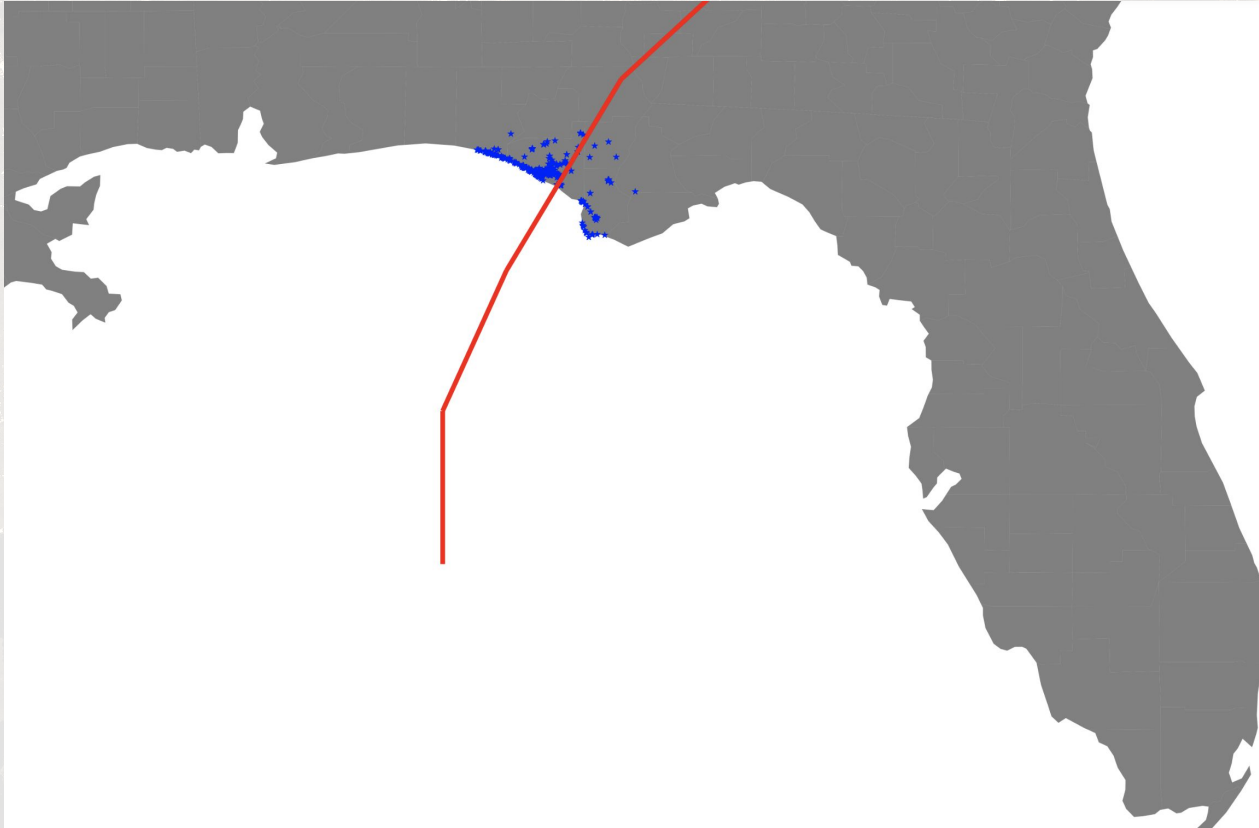
# Mexico Beach after Hurricane Michael



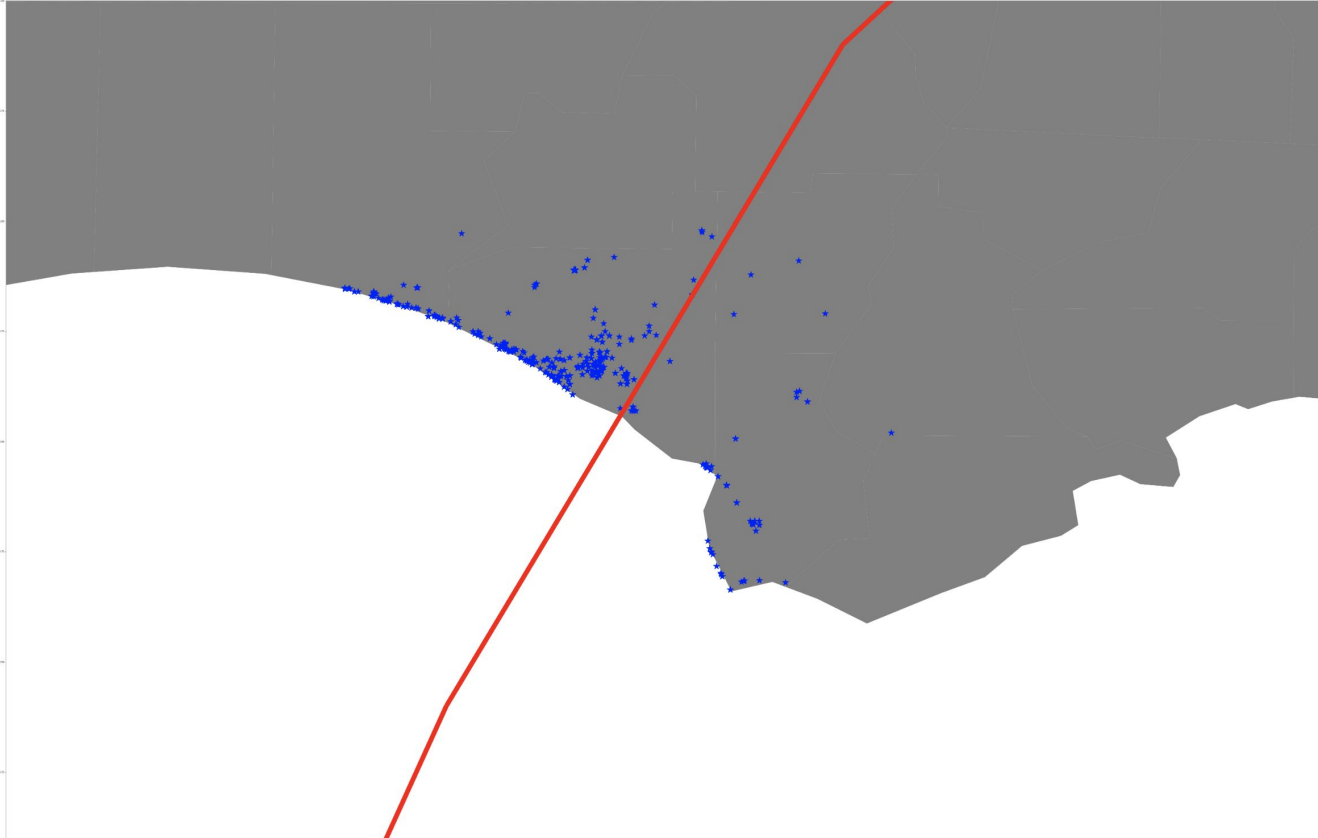Category ● 4 ● 3 ● 2 ● 1 ● Tropical storm     Forecasted path --

Tropical storm-force winds

Hurricane force winds

Source: NY Times

Tweets along path of Hurricane Michael

Mexico Beach, Panama City & Beach

Hurricane Michael - Mexico Beach, Fl

Landfall 2pm EDT (18:00 UTC) October 10, 2018

# Findings

- There are clear linguistic differences in urgent messages vs. non-urgent messages.
  - Sometimes, location is the only thing that can distinguish:
  - *"Earthquake in Haiti. Can't join anybody there I'm dying right now God help Haiti I can't believe it can't stop crying..."*
- Models do well at learning text data associated with labeled groups
  - Can be used on unrelated data and mirror similar trends.
- However... this is only reliable when groups are rigorously classified in the training dataset.
  - Our model suffered from noise due to poorly binned "urgent" messages and inconsistent linguistic patterns found on social media.

# Recommendations & Next Steps

- Disaster response organizations should partner with social media to get access to geolocation data during natural disasters.
  - Would allow for expedient classification and contact with survivors.
  - Security issues, opt-in in user settings
- Clustering Algorithms with sentiment analysis and Word2Vec.
  - Message sentiment is obvious to humans, but we want to train machines to pick up on language context more broadly.
  - If built well, would solve the pre-labeling problem.
- Pay for Twitter's Enterprise API

# Q&A

GA | Data Science Immersive | P4