# "Crash and Burn"

● ● ●

Categorizing "disaster" tweets using Unsupervised Learning models
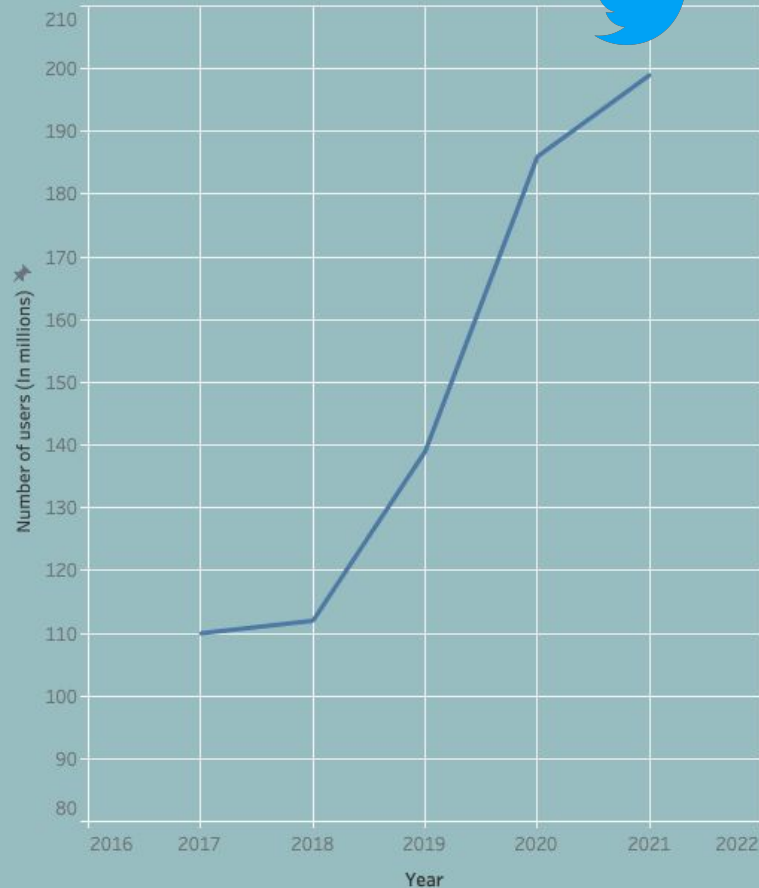
# The Ocean of Twitter

500+ million tweets/day

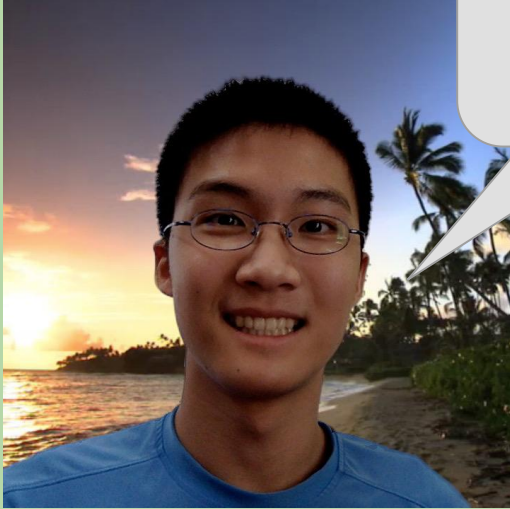200+ million users

Countless categories that
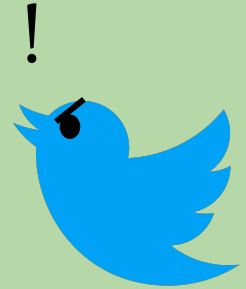can be modeled



Number of Twitter Users

# Objective: Create an unsupervised model to categorize tweets into topics
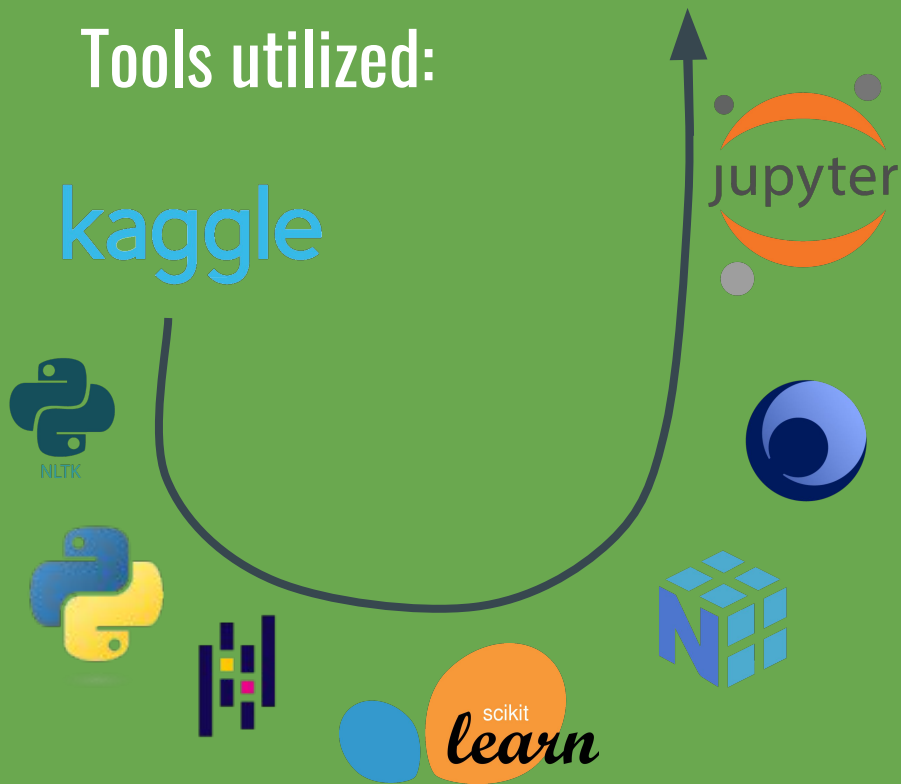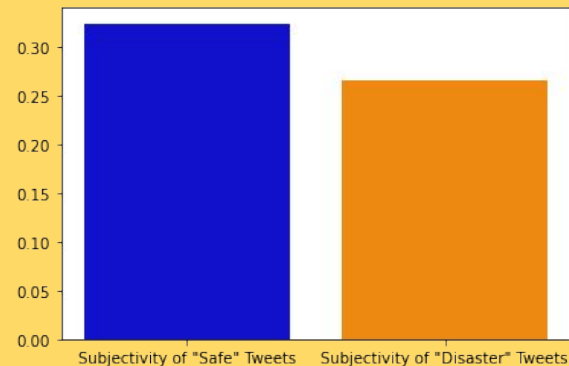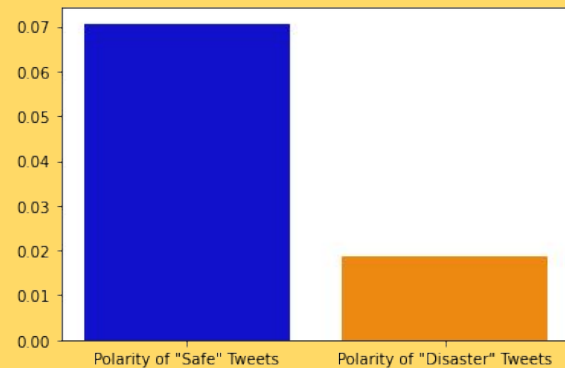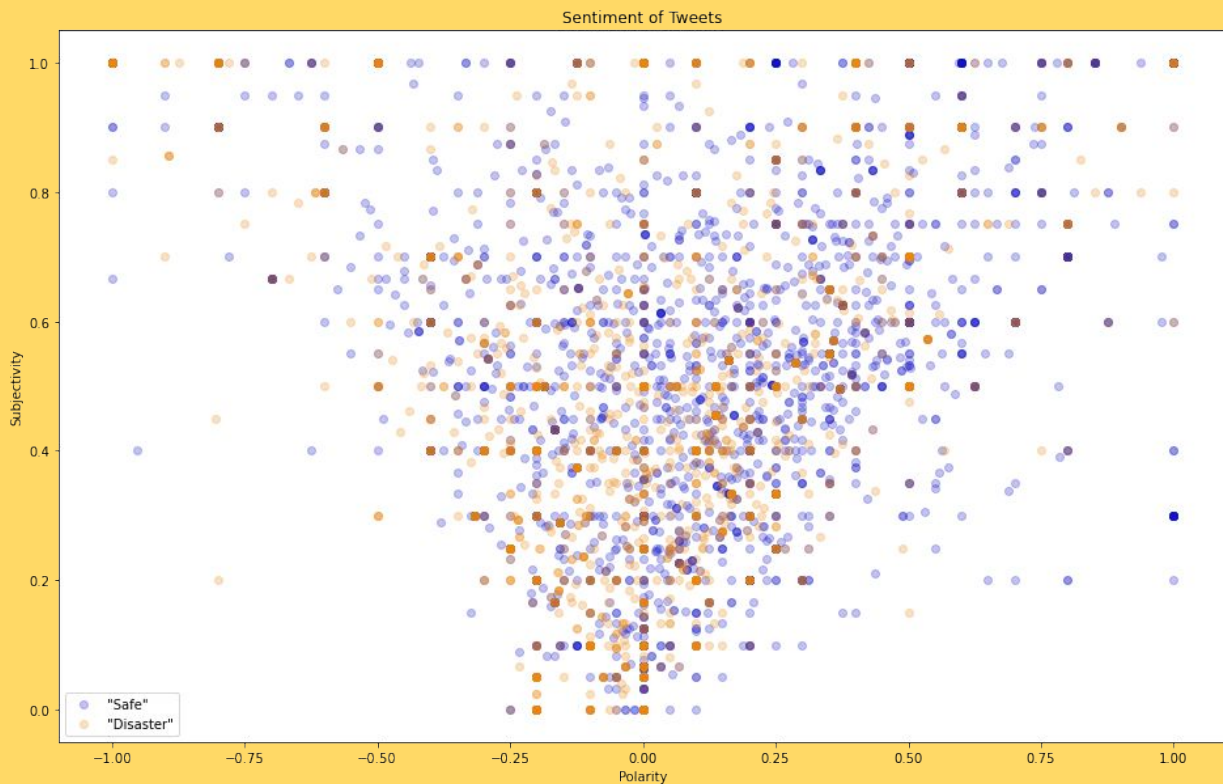
# Data

- 7,613 tweets
- 4,342 "safe"
- 3,271 "disaster"

# Tools utilized:

# Preliminary Findings

# Preprocessing Techniques

Default Tokenizer

Lemmatization

Cleaning non-Unicode characters

Using bigrams as tokens



NLTK

# Model

LDA Topic Modeling

K-Means Clustering

Max Frequency = 7.5%

Min Frequency = 0.5%

8 topics

# Topics - most common words

| | |
|---|---|
| 'Accidents': help, top, near, keep, collide, person | 112 |
| 'Fire': burn, see, wild, wild fire, school, car, high | 235 |
| 'Injury': wound, blood, movie, fatal, building | 208   ("movie"? Mislabeled tweets contaminating) |
| 'Police': police, officer, home, collapse, area, line | 249 |
| 'Storms/Flooding': wreckage, weapon, injury, hiroshima, nuclear, damage | 198   (Japan-related overlap) |
| 'Terrorism': man, attack, bomb, tragedy, city, fall | 302 |
| 'War': war, emergency, kill, talk, loud, service | 131 |
| 'Misc': Everything else | 432 |
| | |
| Total: | 1866 |

# Using the model

The model can be used to categorize any text into the eight main clusters.

It can be easily accessed from a secondary notebook in the GitHub repository.

# Future Work

- More data, without unrelated topics
- Seed topic weights to reduce topic overlap
- Explore other models

"going to redo my nails and watch behind the scenes of desolation of smaug ayyy"

"Next May I'll be free...from school from obligations like family.... Best of all that damn curfew..."

"OMG?? Didnt expect Drag Me Down to be the first song Pandora played OMG I SCREAMED SO LOUD My coworker is scared"

"@Drothvader @CM_Nevalistis you can keep this please!!!!! Arachys [2 Pieces] - Now deals 4000% weapon damage (up from 2500%)"

Questions?