

NLP Text Classification & Sentiment Analysis on Disaster Tweets



Nicole Hong

LHL Demo Day: April 28, 2022



GettingStarted Prediction Competition



Kaggle · 837 teams · Ongoing

Natural Language Processing with Disaster Tweets

Predict which Tweets are about real disasters and which ones are not

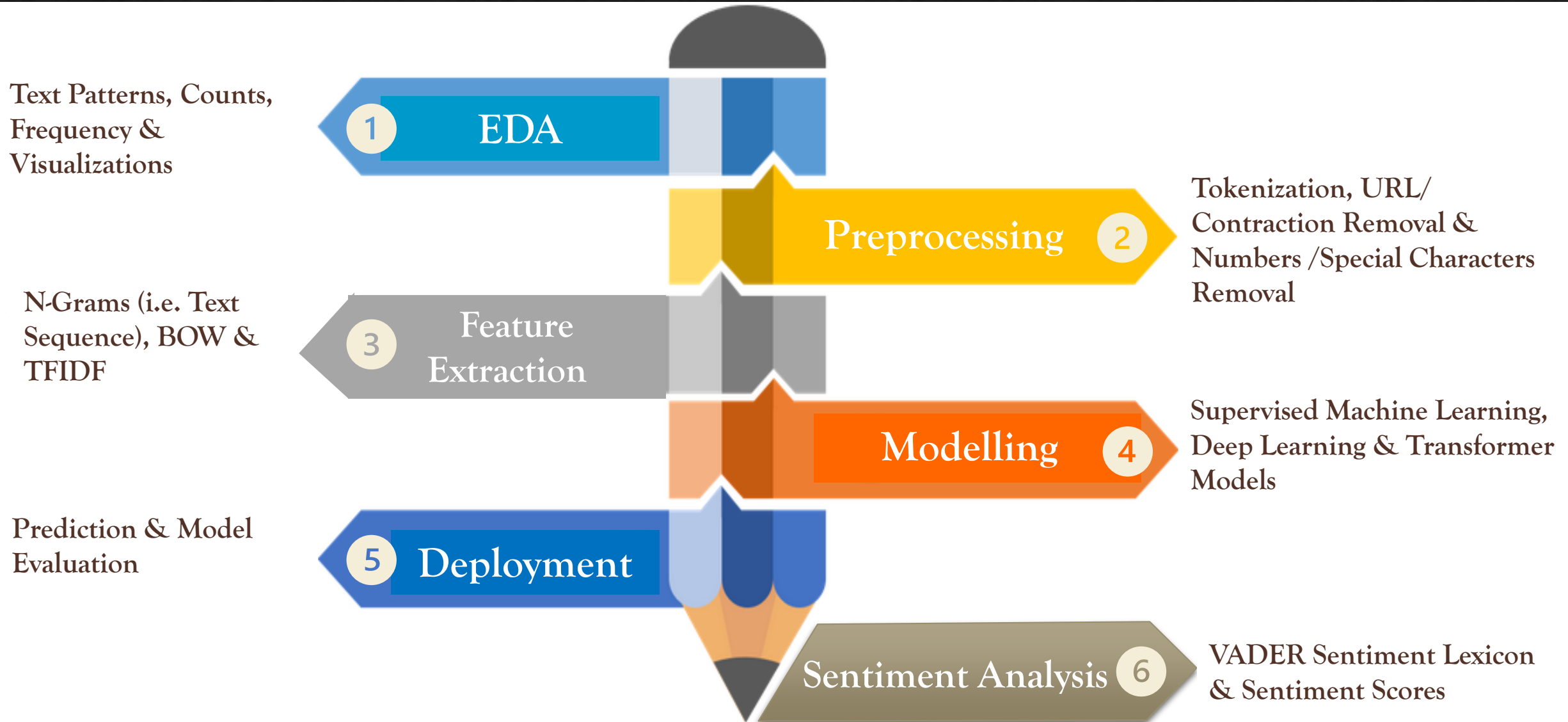
Why this Kaggle?

- Twitter as an important communication channel in times of emergency
- Smartphones enable announcement of an emergency in real-time
- More agencies interested in Twitter & Emergency Response Management

Objective

- Identify disaster related tweets with high accuracy (Classification)
- Determine level of crisis for better response to disasters (Sentiment Analysis)

Project Workflow

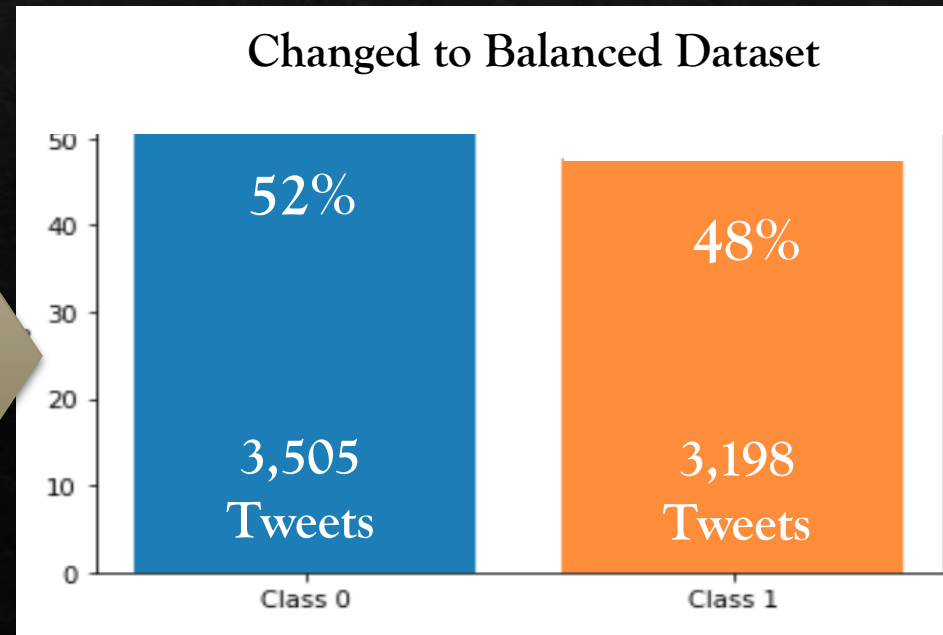
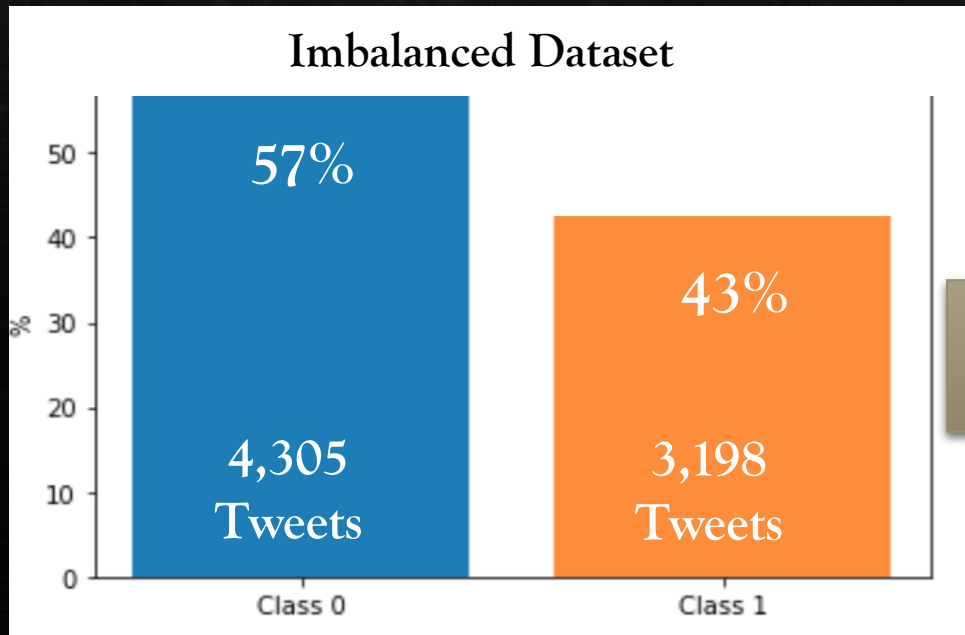


About Kaggle Data

- Small Text Data:
 - Train Data = 7,613 entries
 - Test Data = 3,263 entries

- Variables with Keywords, Locations, Text and Target
- 28% of Missing Data in Keywords & Location
- Tweets in August 2015
- Not designed for Sentiment Analysis

Target Class Balance – Training Dataset



Class 0

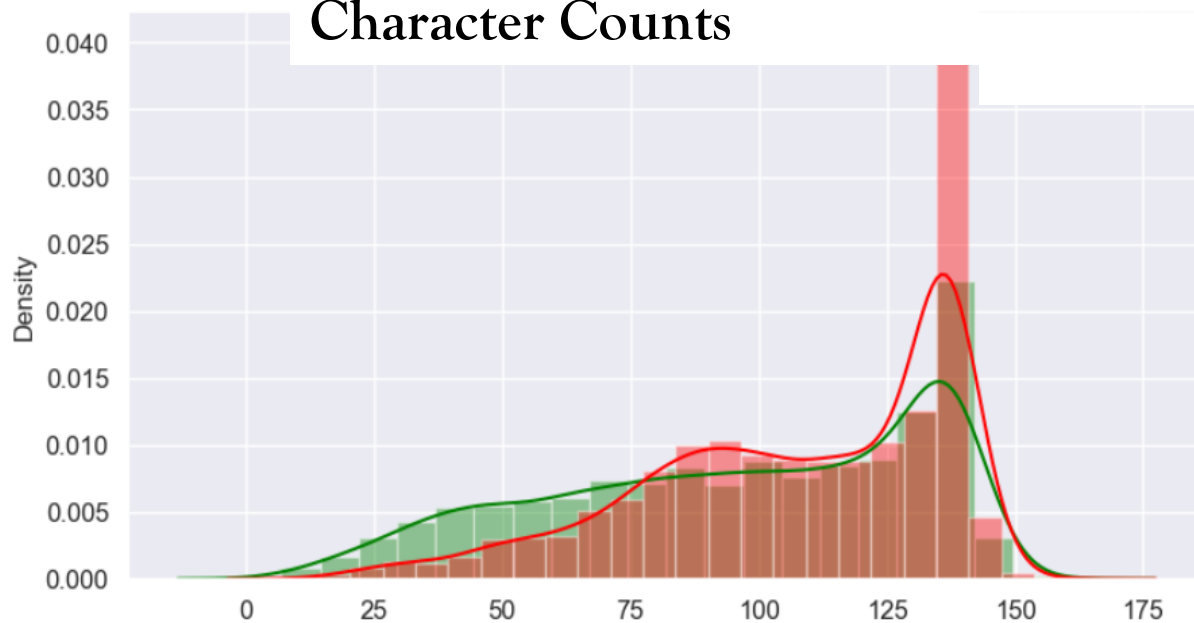
Non-disaster
related Tweets

Class 1

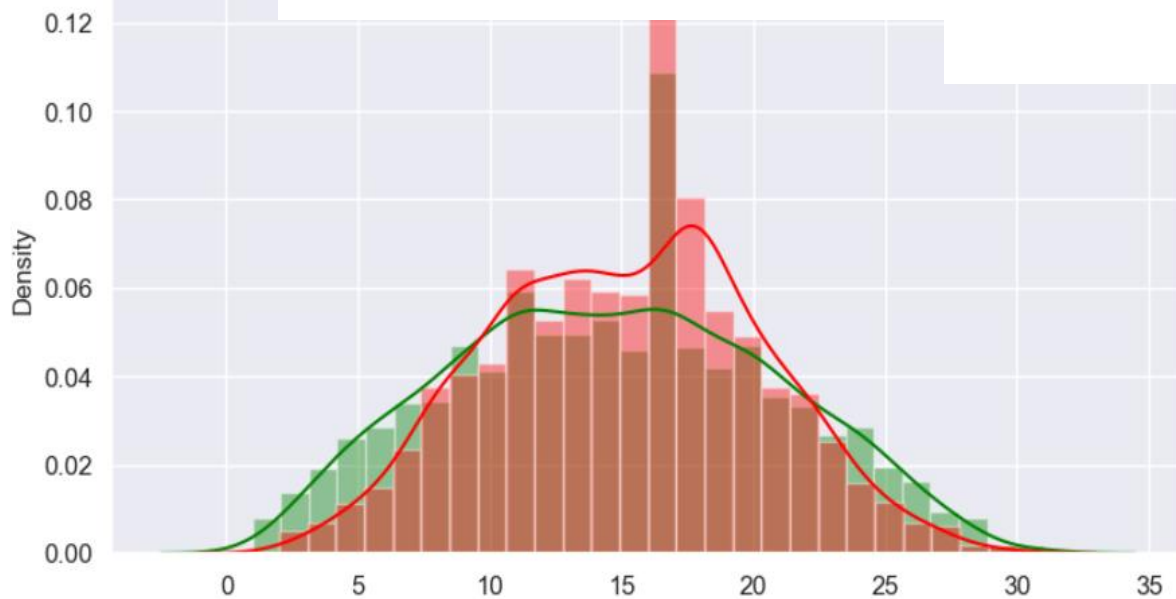
Disaster related
Tweets

EDA: Distribution Plots

Character Counts



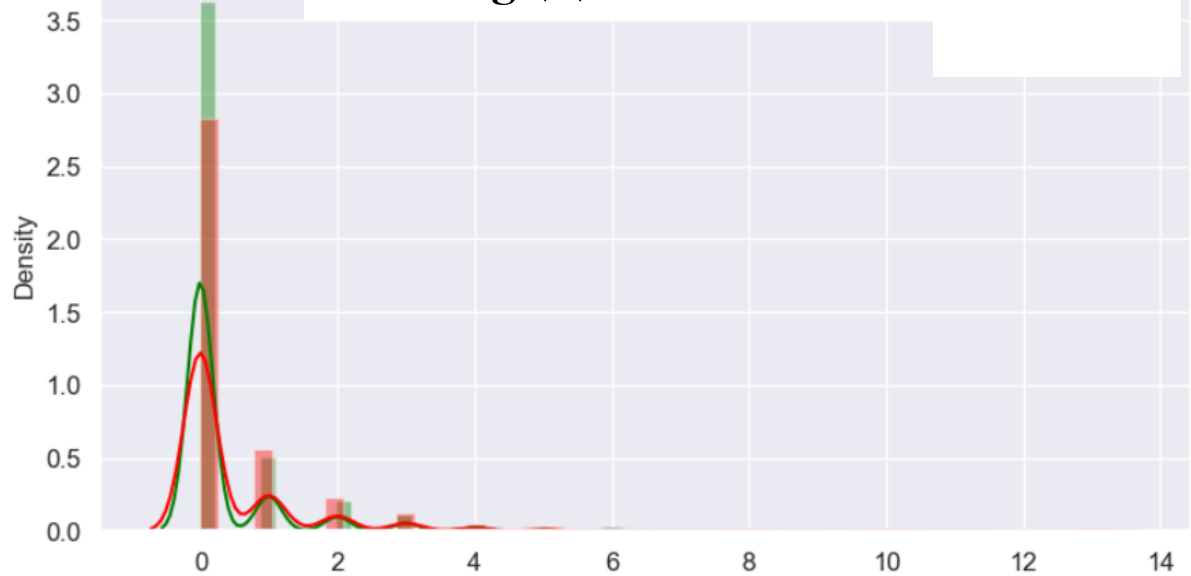
Word Counts



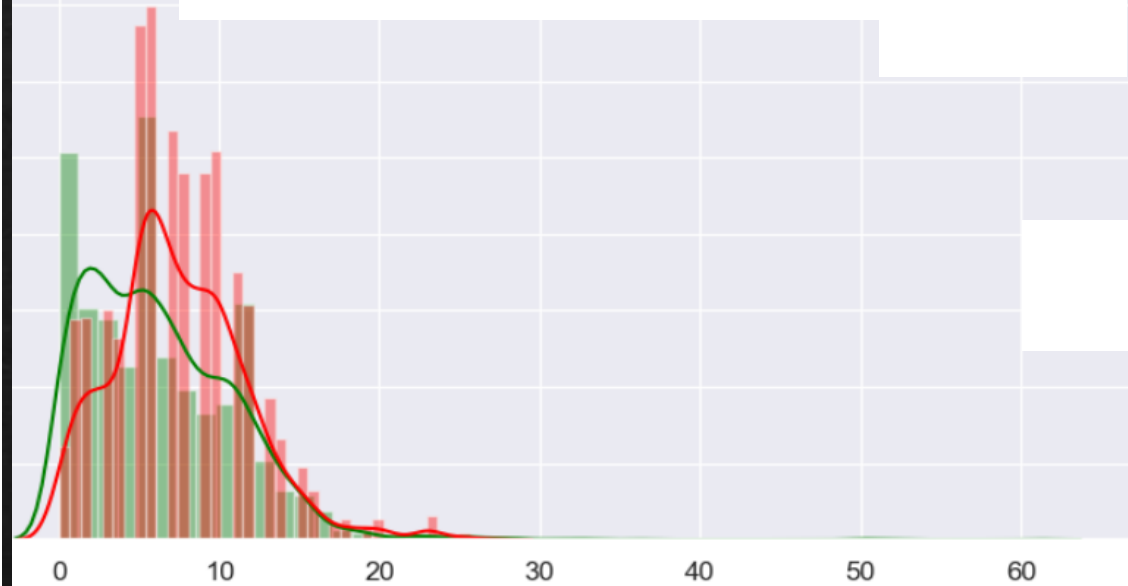
Stop Word Counts



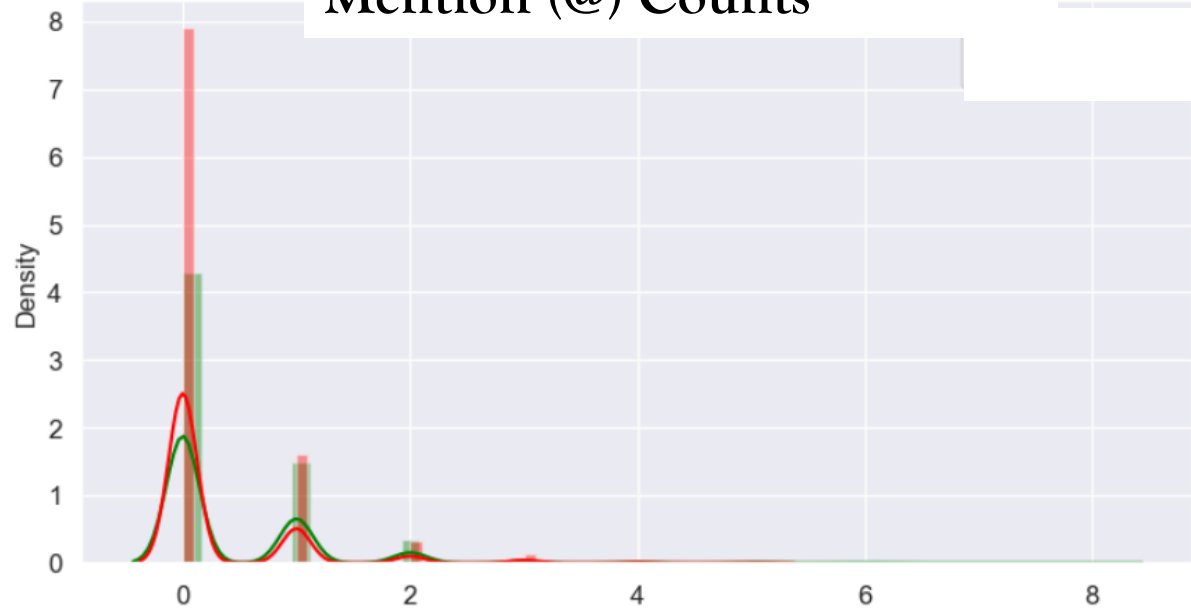
Hashtag (#) Counts



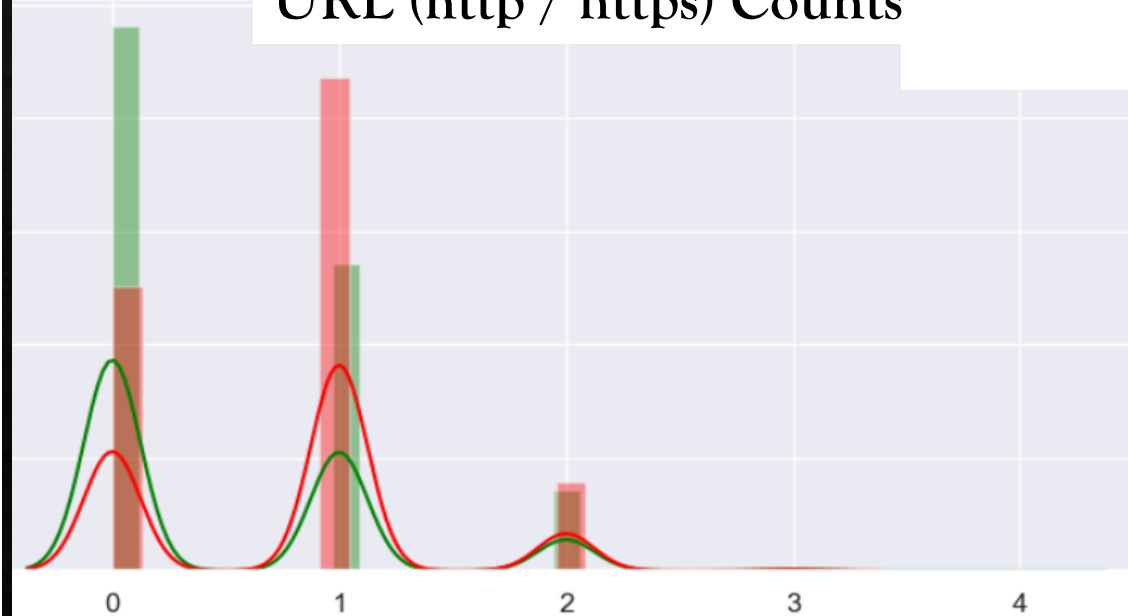
Punctuation Counts



Mention (@) Counts

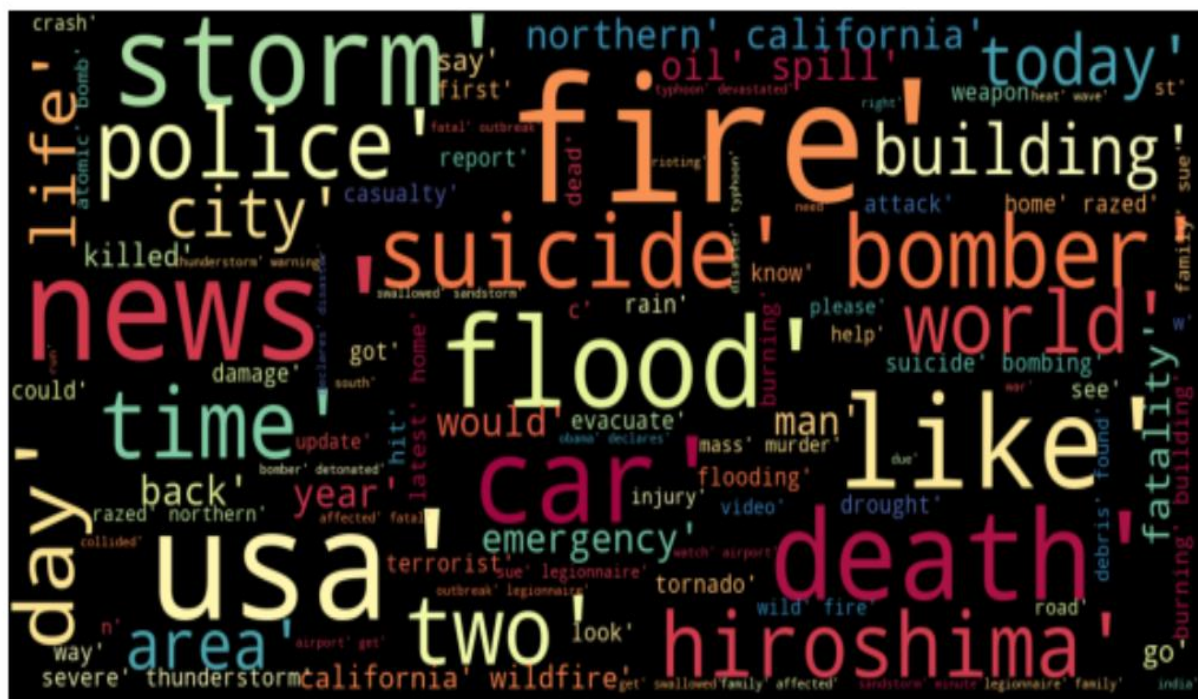


URL (http / https) Counts



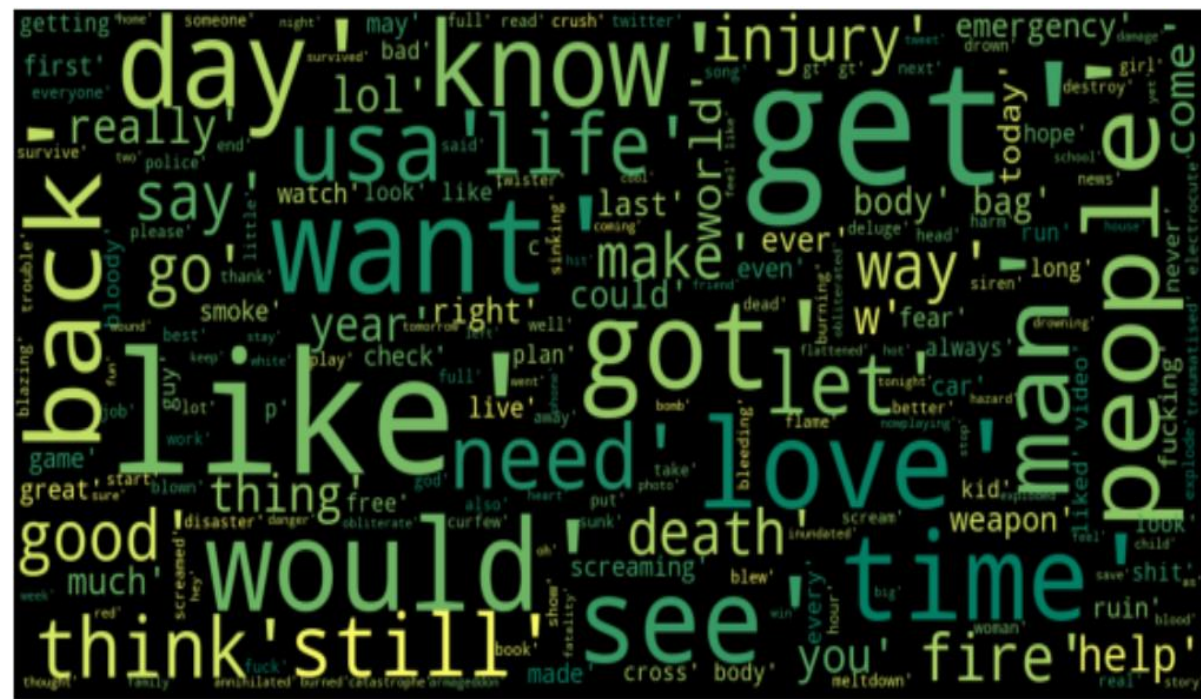
WordCloud: Commonly Used in Tweets

Disaster Tweets | Target = 1



Topic	Number of Tweets
fire	263
news	146
police	108
usa	108
suicide	106

Non-Disaster Tweets | Target = 0

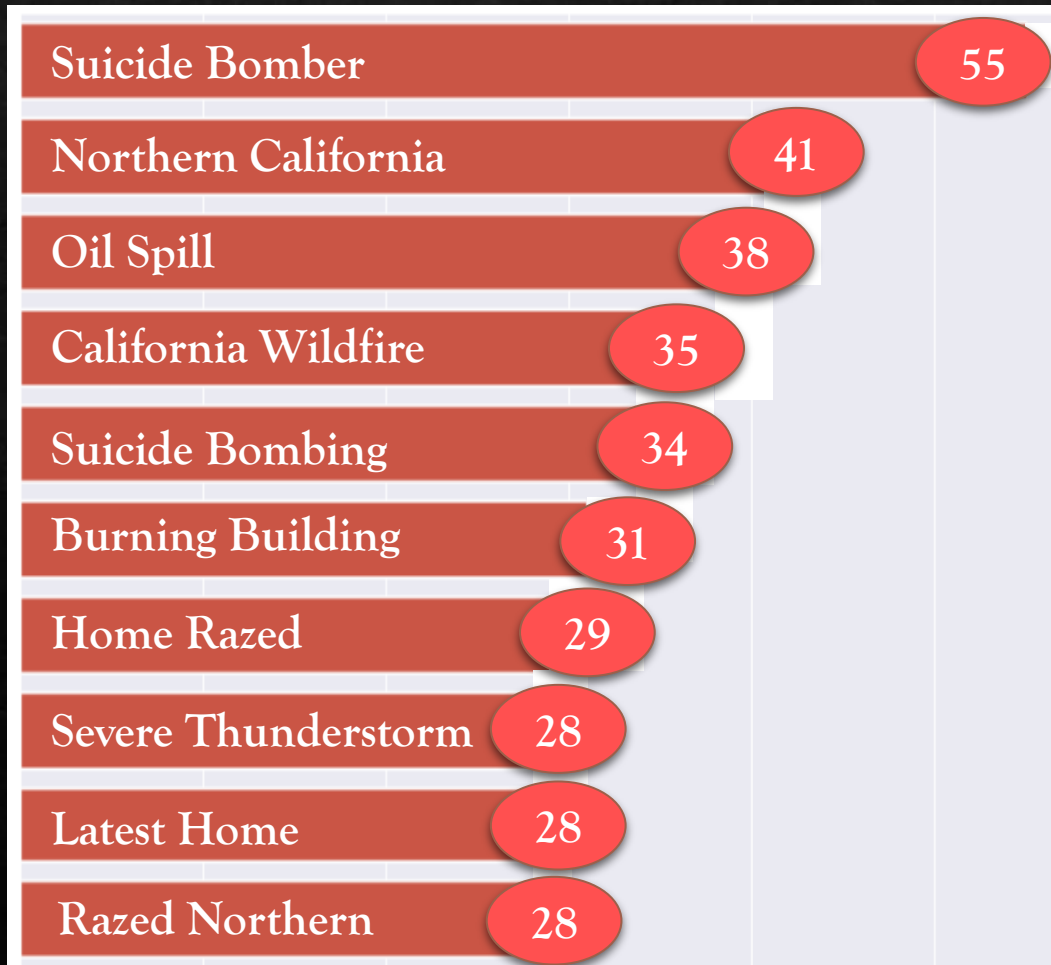


A horizontal bar chart with five bars of varying lengths, colored in a dark blue. The bars are labeled on the left with the verbs 'like', 'get', 'want', 'time', and 'would'. The corresponding numerical values are displayed at the end of each bar: 206 for 'like', 155 for 'get', 89 for 'want', 83 for 'time', and 82 for 'would'. The chart is set against a light gray background with vertical grid lines.

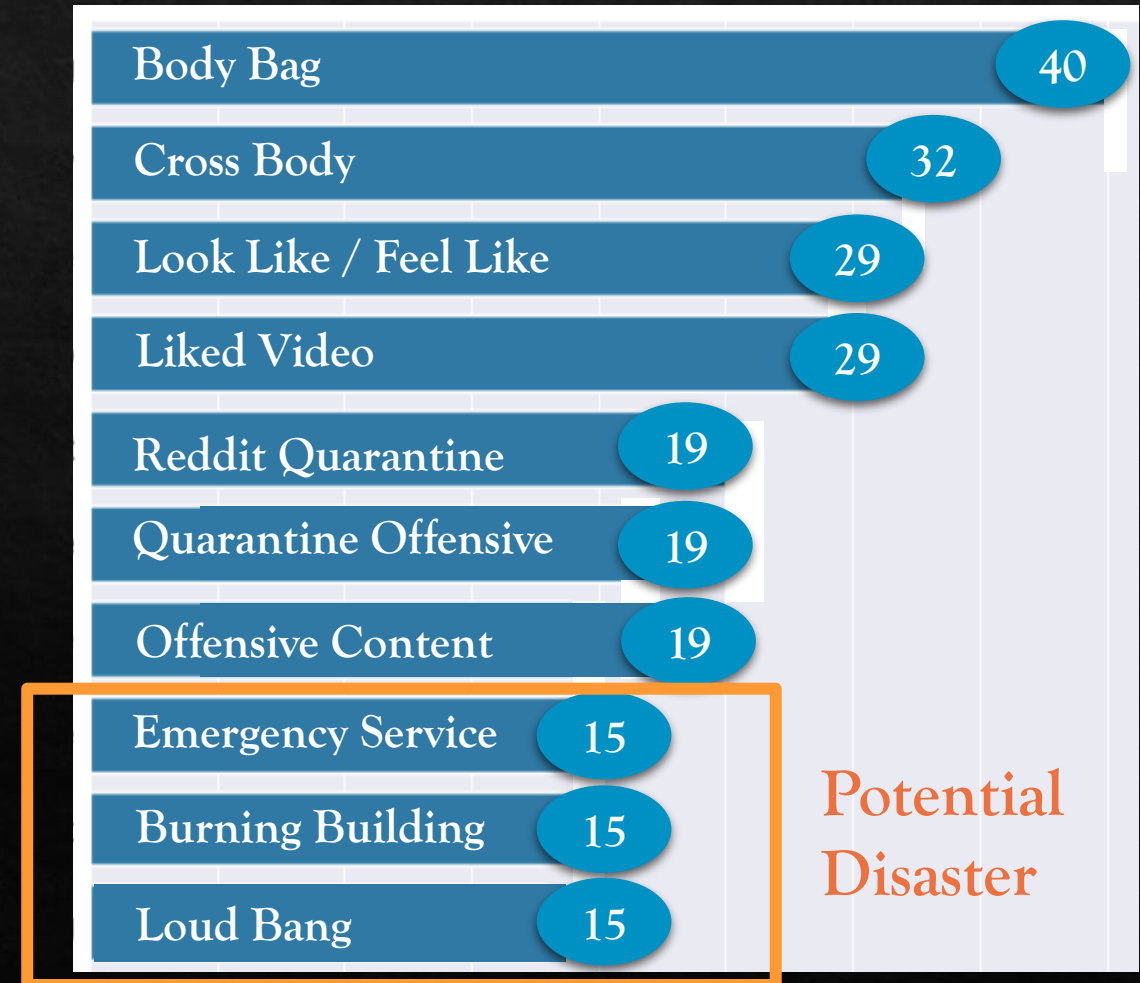
Verb	Count
like	206
get	155
want	89
time	83
would	82

Representation: Text Sequence (N-Grams)

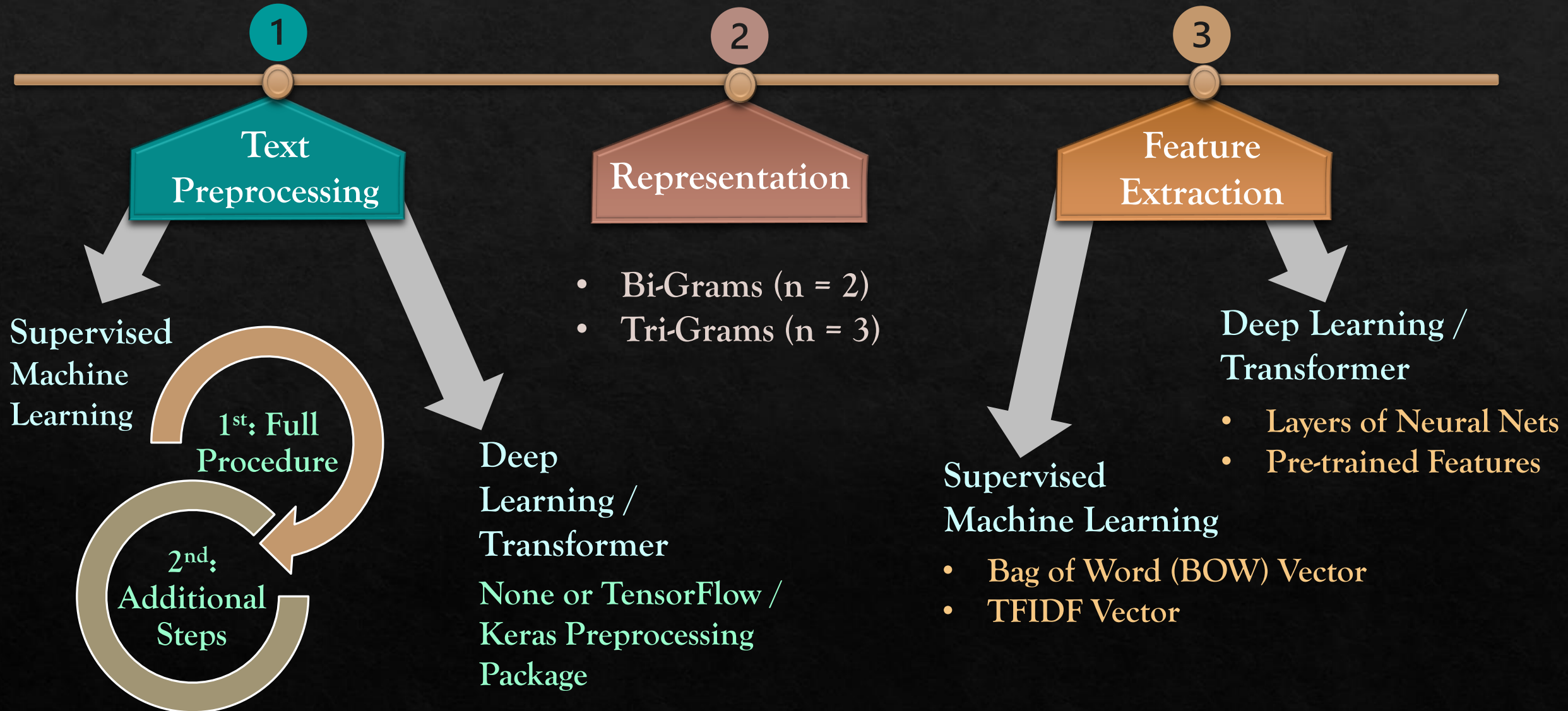
Disaster Tweets | Target = 1



Non-disaster Tweets | Target = 0



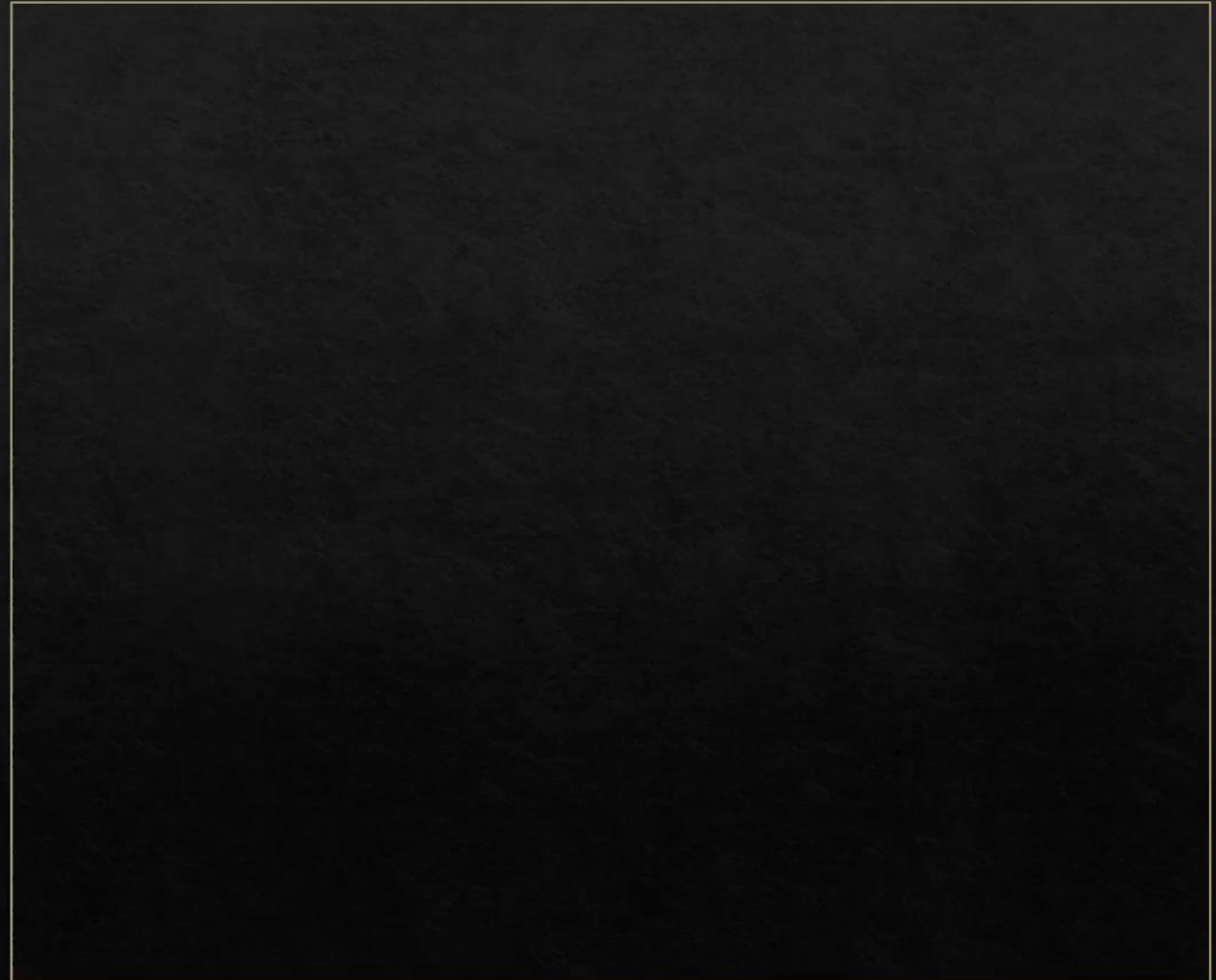
Basic Steps for Text Classification



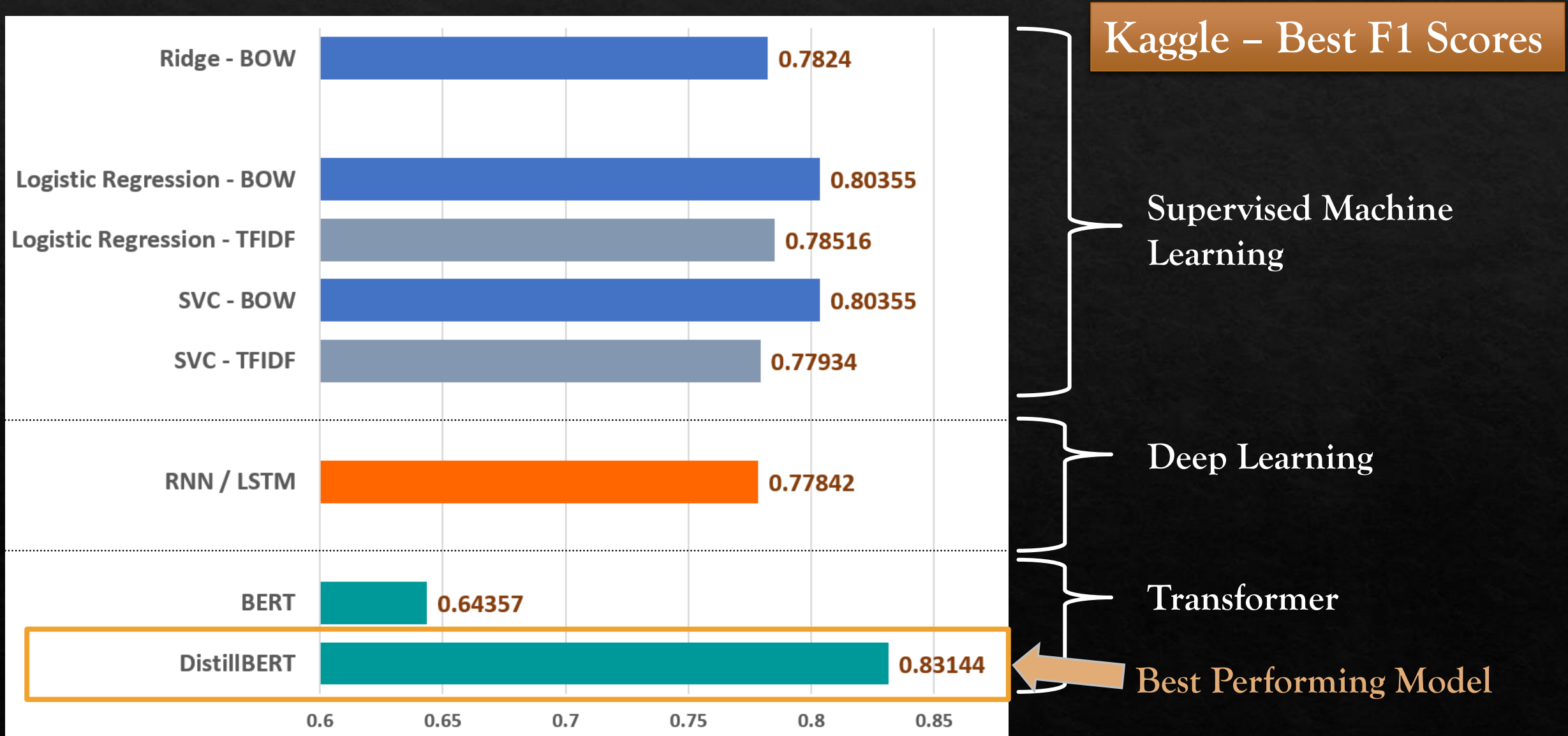
Text Preprocessing

1st Preprocessing:

- Remove Duplicate Tweets
- Replace Abbreviated/ Special Text to Standard Text
- Remove URL
- Remove Contractions
- Tokenize
- Remove Numbers & Unicode
- Remove Punctuation
- Normalize
- Clean Repeated Characters in Text
- Remove Stop Words
- Lemmatize
- Remove Contractions
- Join Tokenized Words to Sentence



Machine Learning – Classification Performance

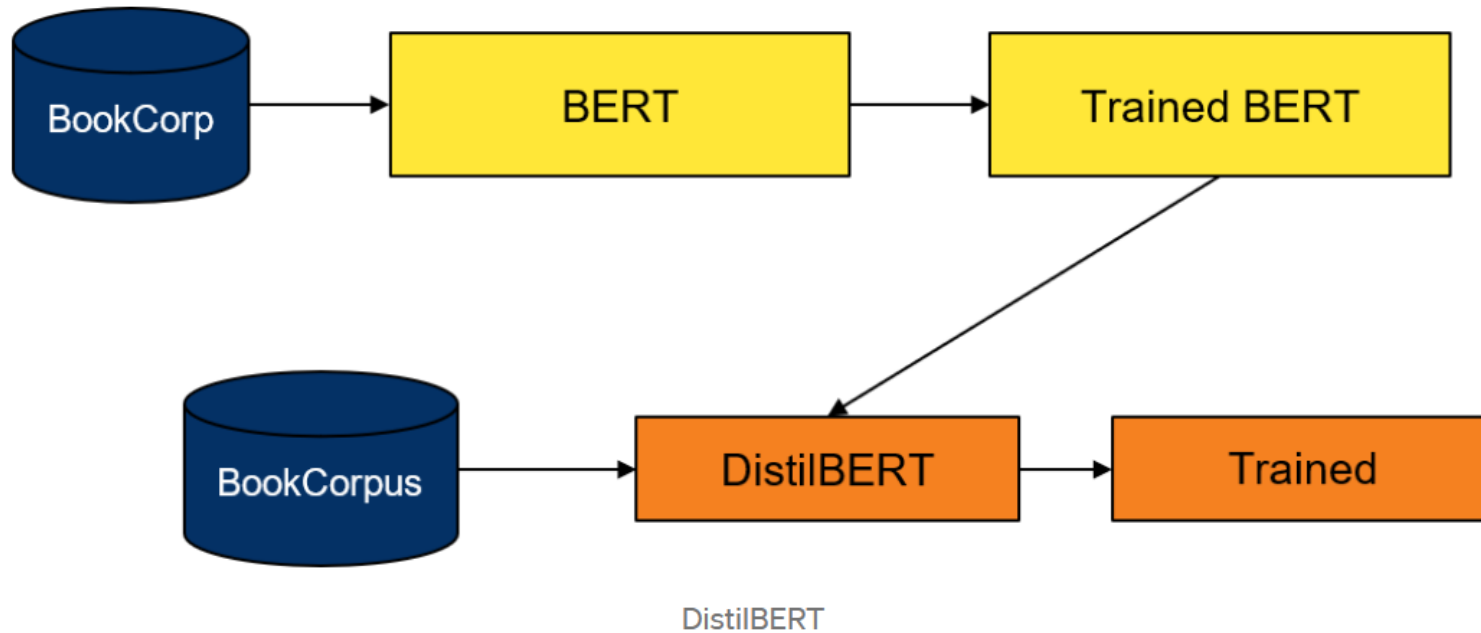


Transformer - DistilBERT



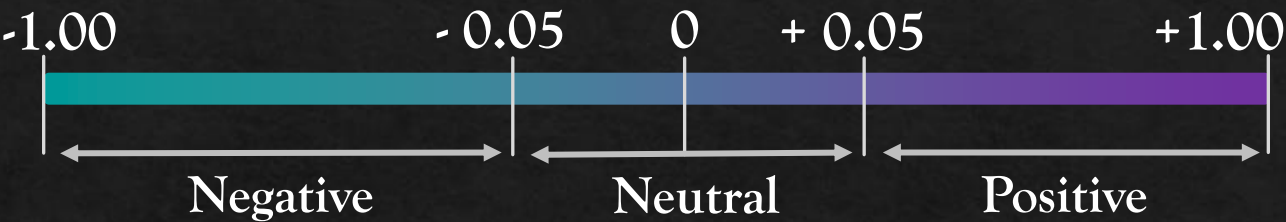
Hugging Face

Smaller, lighter and faster version of BERT developed and open-sourced by the team at HuggingFace

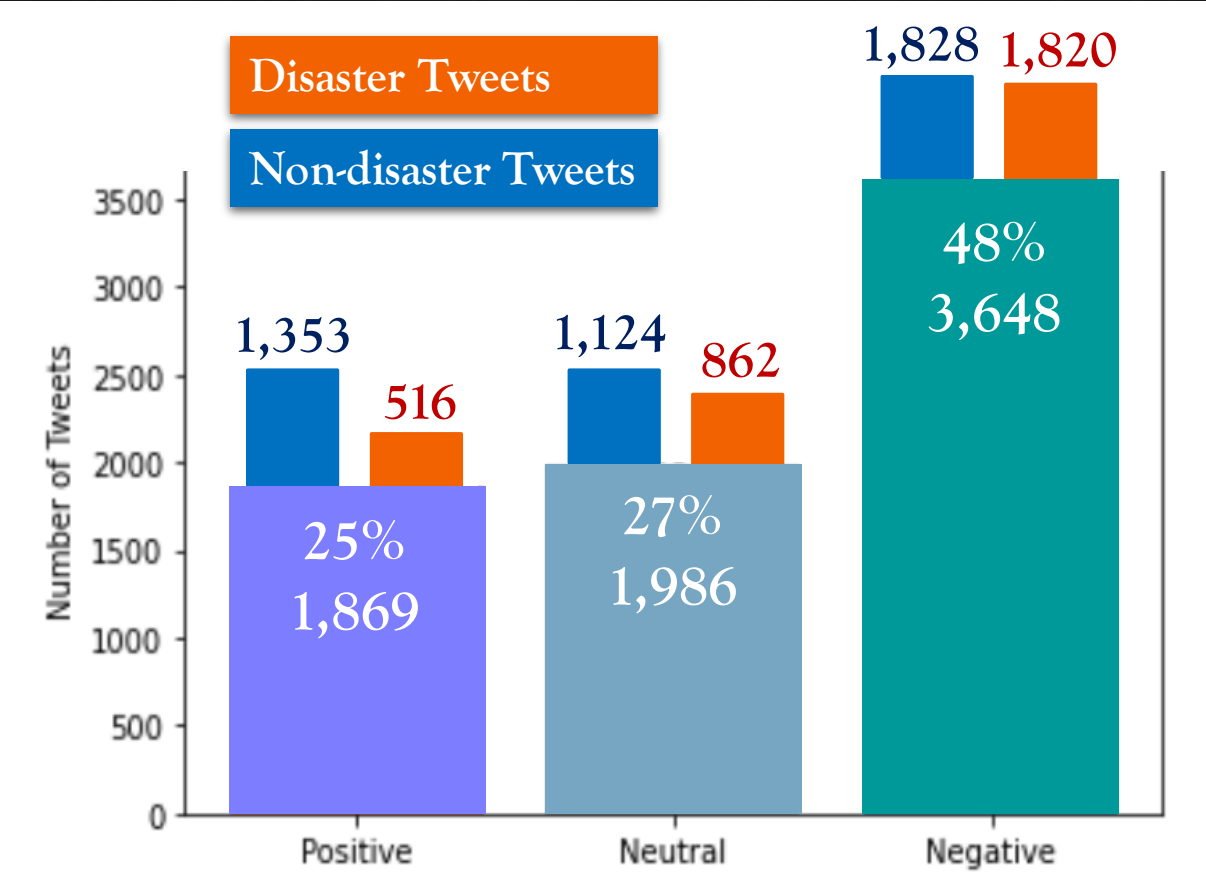


- 40% less parameters than BERT-base-uncased, runs 60% faster, preserving over 95% of BERT's performances
- Implementation on TensorFlow/Keras, Trained DistilBERT used to generate sentence embedding
- Basic NN Architecture (with Dense and Dropout layers)

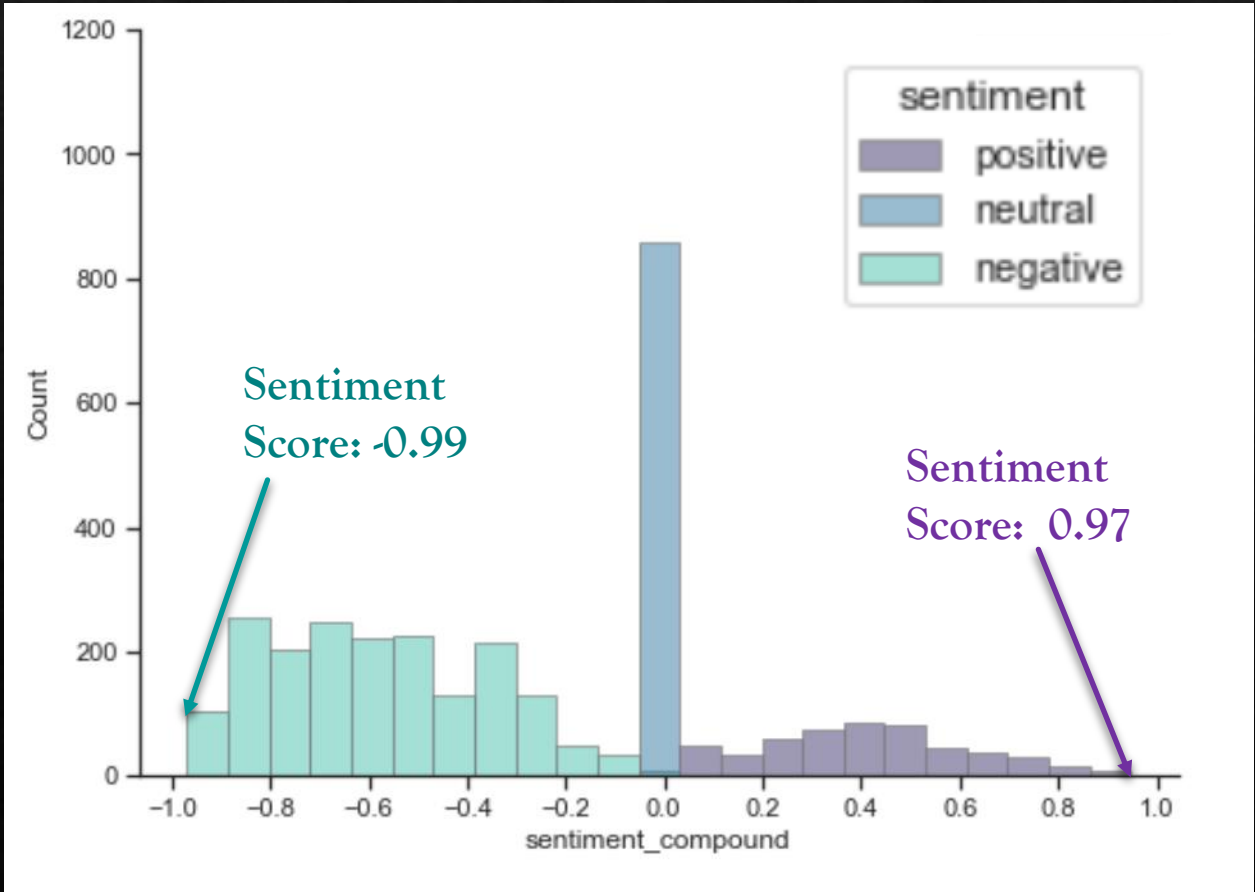
Sentiment Analysis



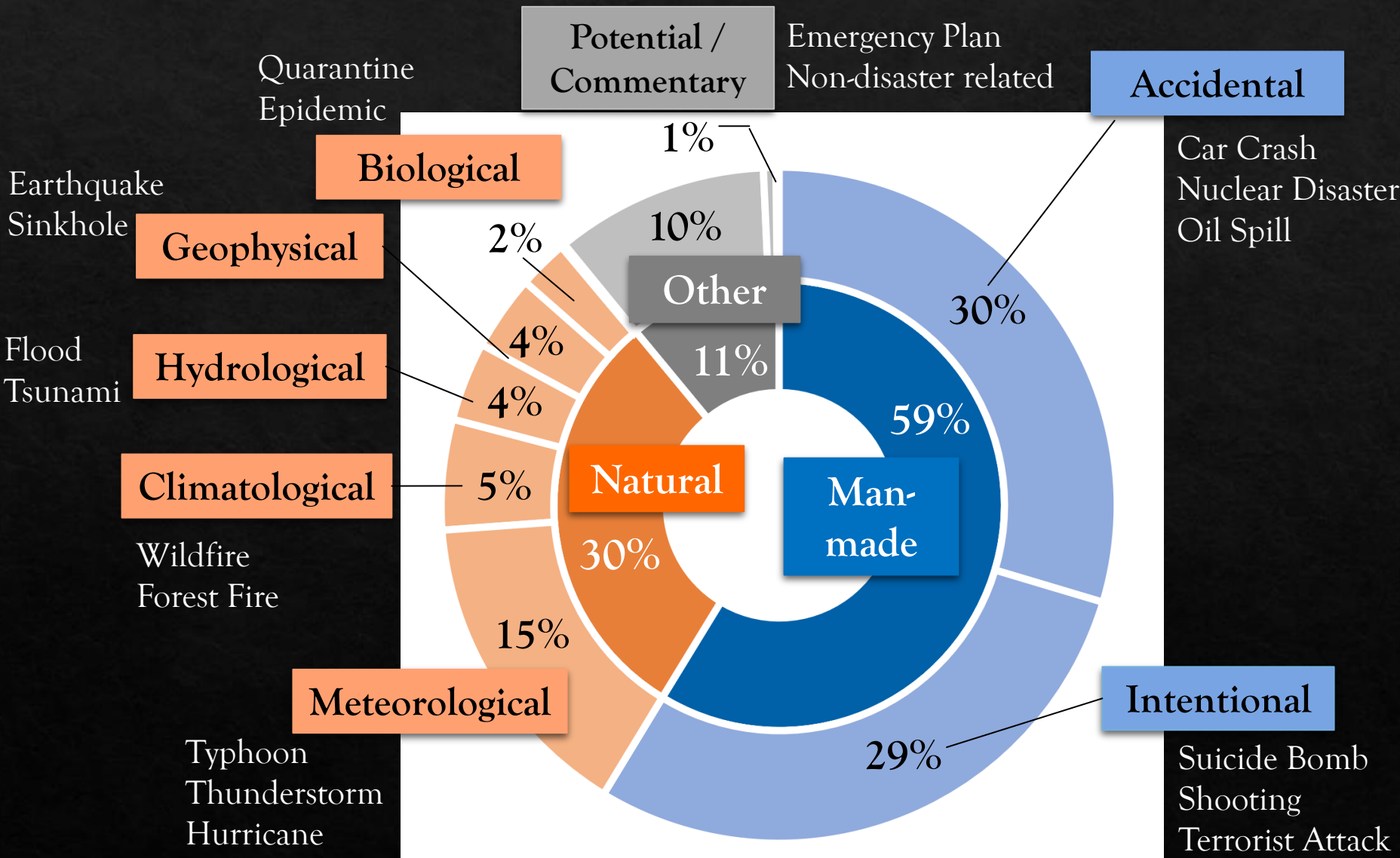
Counts of Sentiment Classes in All Tweets



Sentiment Score Distribution – Disaster Tweets



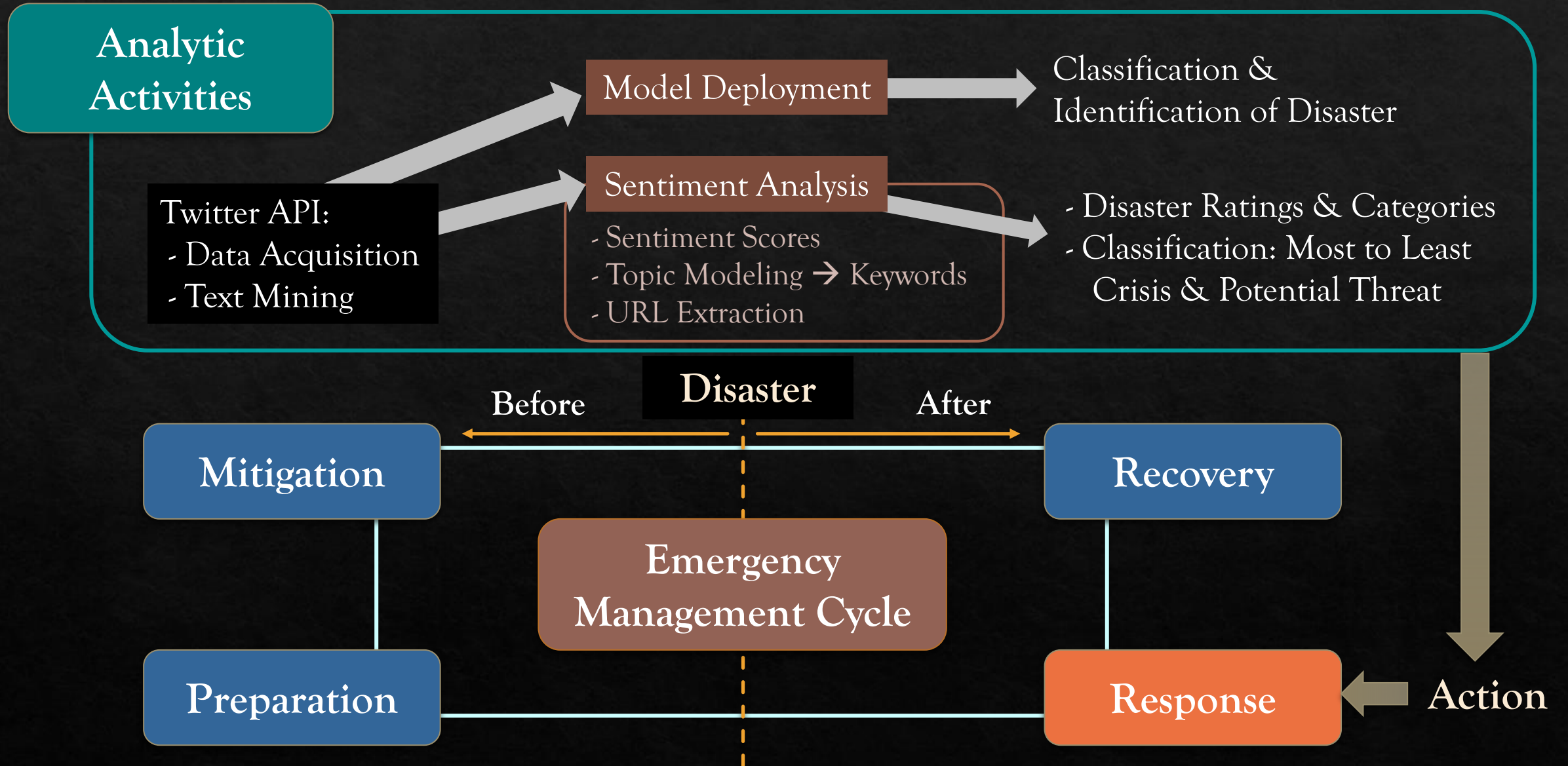
Disaster Tweets

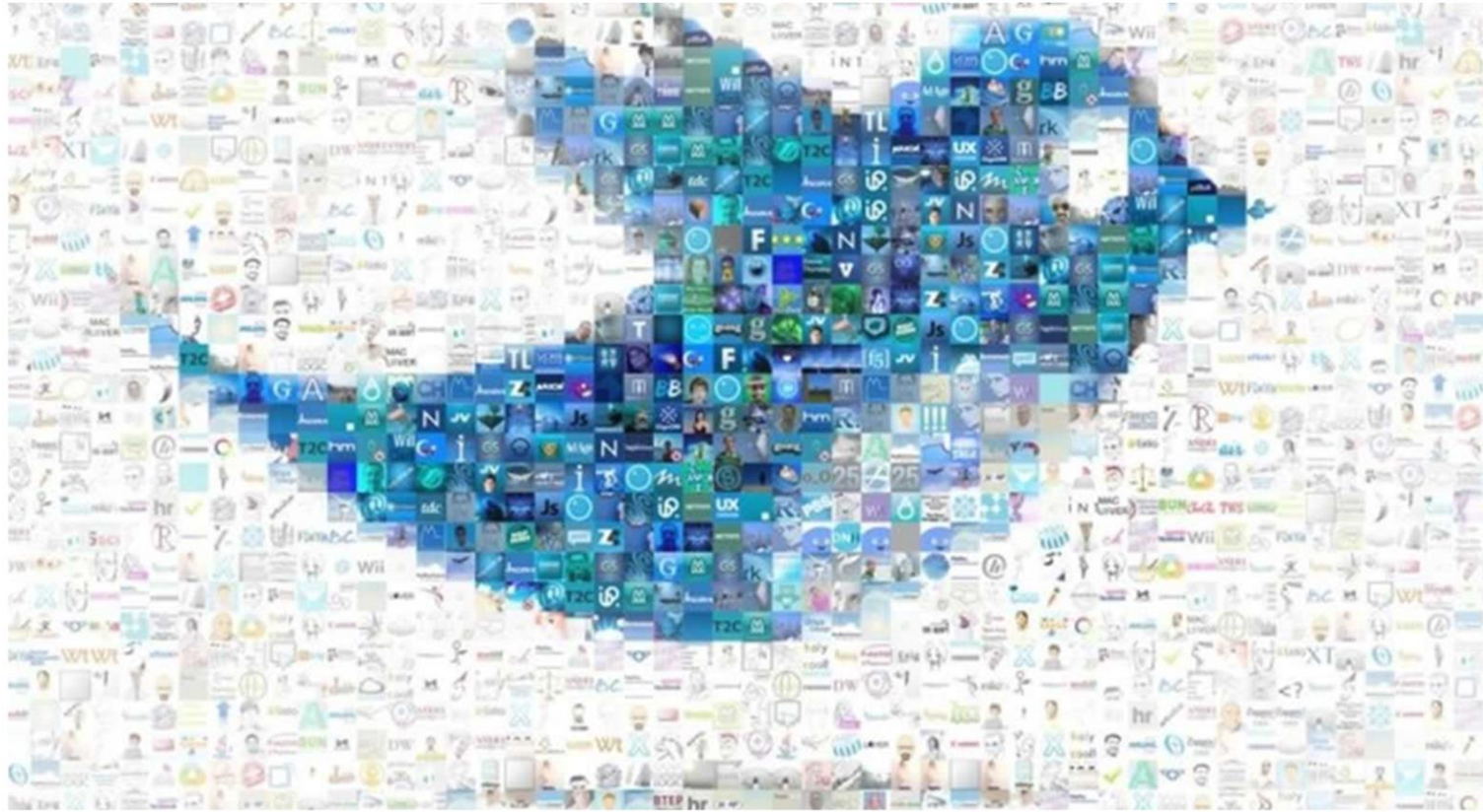


Disaster Class & Categories

Class	Sub-class
Natural	Meteorological Hydrological Geophysical Climatological Biological
Man-made	Intentional Accidental
Other	Non-disaster related

Emergency Response Management





Thank You!

✉ n.hong@queensu.ca

🔗 www.linkedin.com/in/nicolehhong

🔗 <https://github.com/Nicole-Hong>