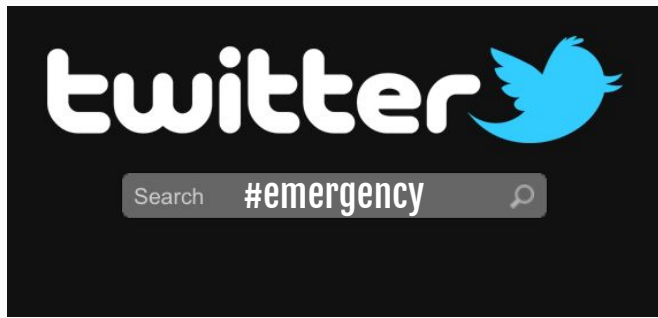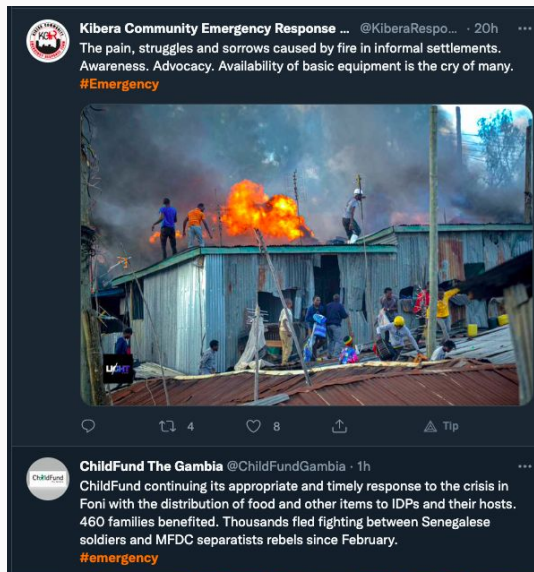# Disaster Tweet Prediction using NLP
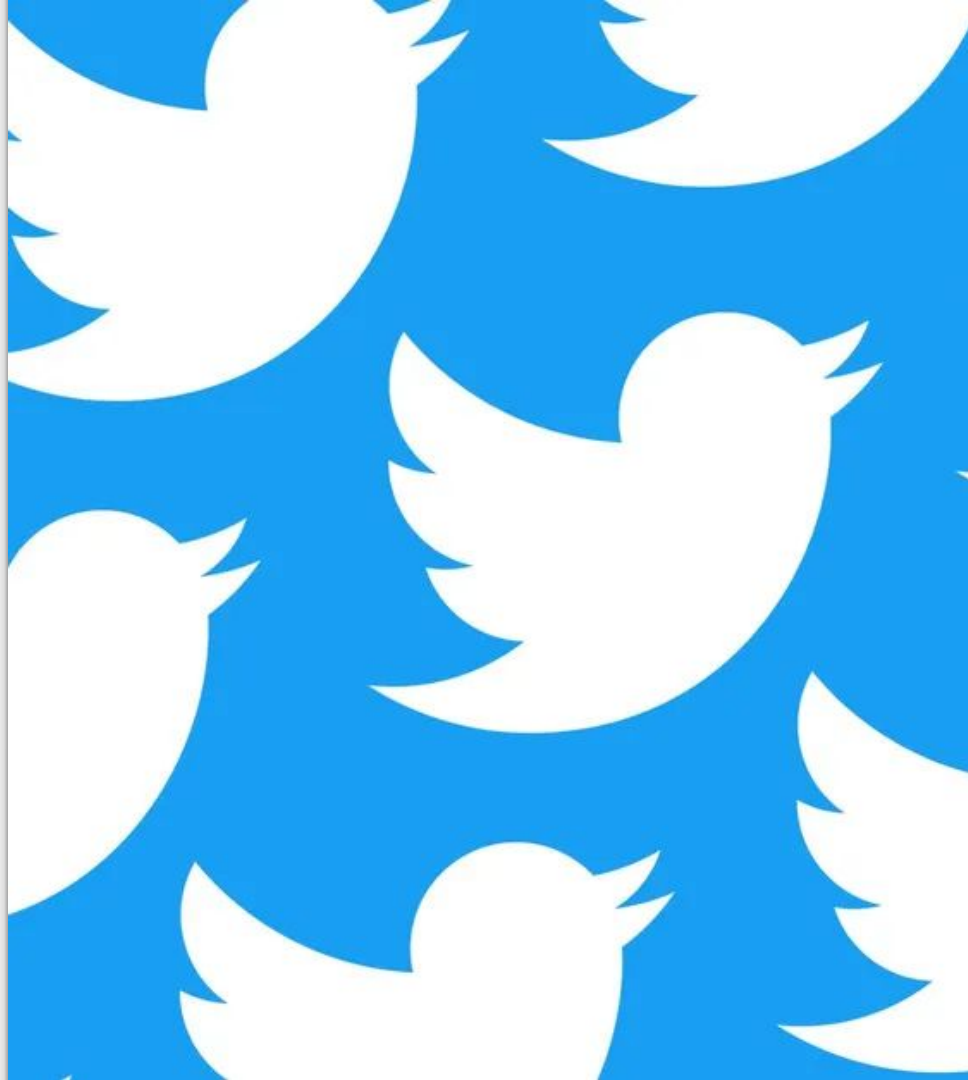
By Tyusha Sarawagi & Sarah Wright

We have seen in recent years how powerful Twitter can be as both a resource for finding urgent information and as a tool for communicating useless information about sales and petty gossip.

Hashtags can be useful tools for sifting through the nonsense, however, oftentimes they are misused
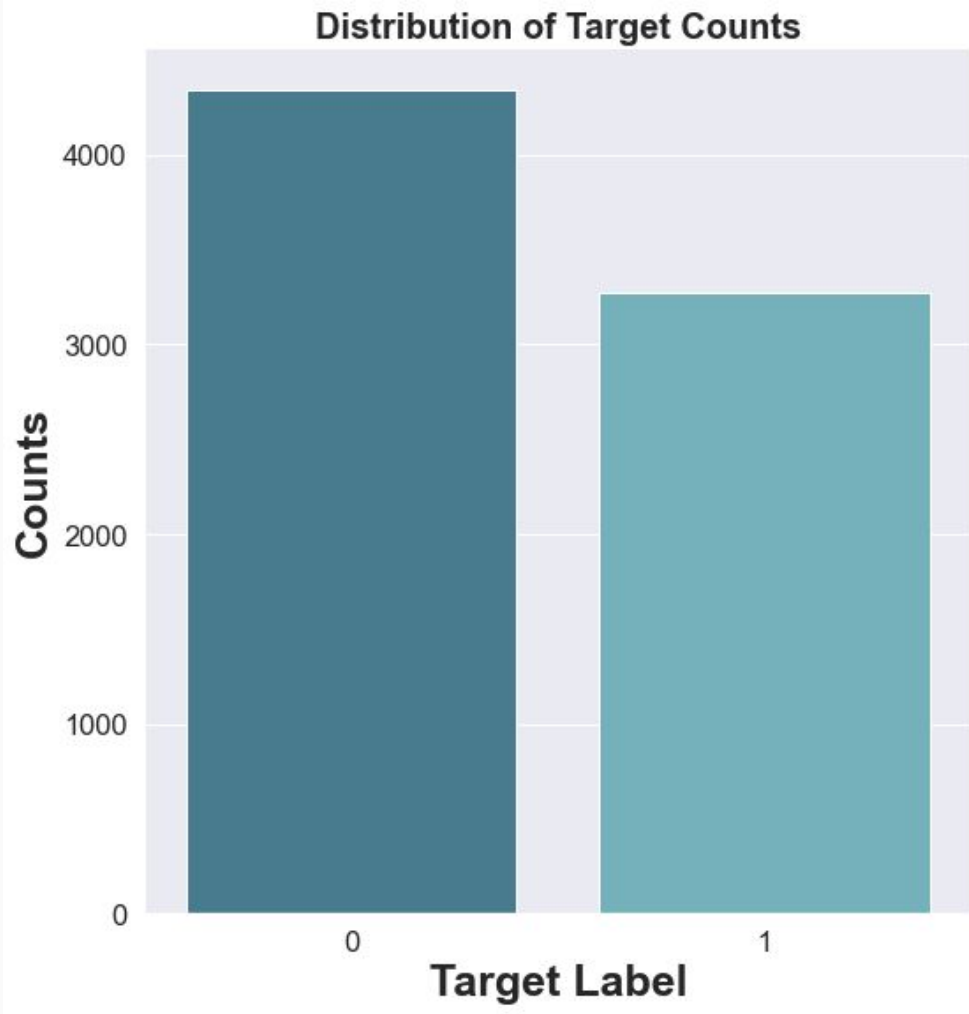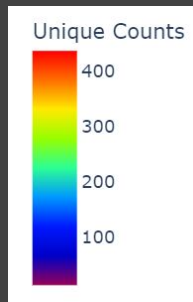
# Disaster Tweet Prediction

# What do we know?

❖ 7613 tweets, each with its twitter-id, location, a keyword chosen from the tweet, and a target indicating whether the tweet is a disaster tweet or not in the training set.
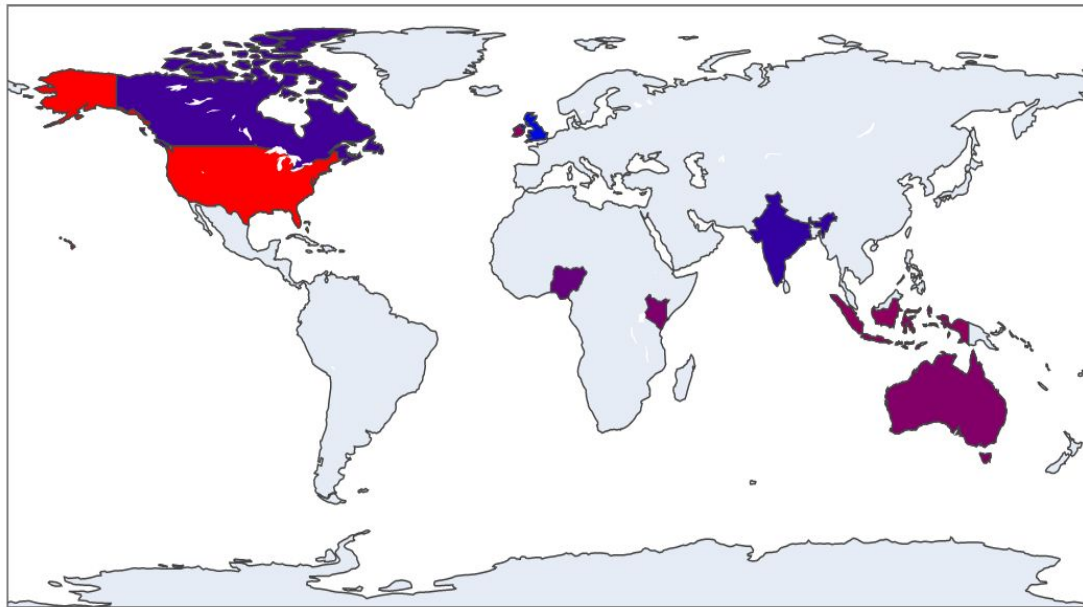


twitter

# What do we know?

- ❖ There are, 4521 unique locations in the dataset

- ❖ Top ten countries from the dataset



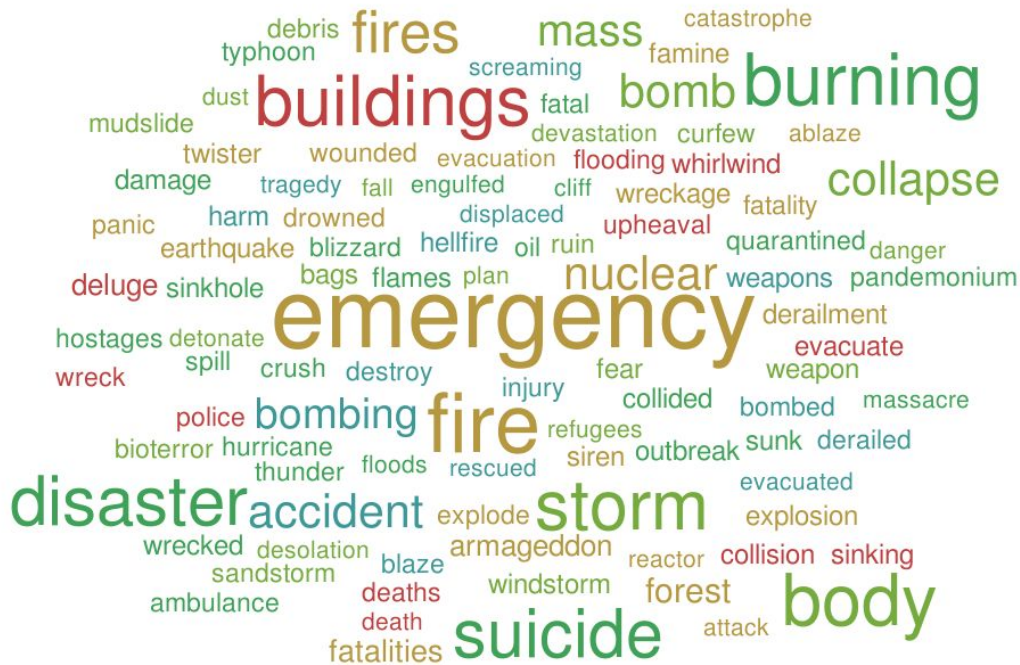Tweet Locations of Top 10 Countries

Unique Counts
400
300
200
100

twitter

# What do we know?

❖ The words "emergency", "disaster", "fire", "storm", "suicide", "burning" occurred more frequently in the text data as compared to others.



twitter

# What do we know?

❖ Non-disastrous tweets (only).
➢ Most Frequent Words: "body" and "emergency".

❖ Disastrous tweets (only)
➢ Most Frequent Words: "suicide", "emergency". "Bombing", "storm", "fire"

➢ There is some overlap, such as with the word "emergency".

Which messages on Twitter can be **classified as Natural Disasters,** and how can we do that?
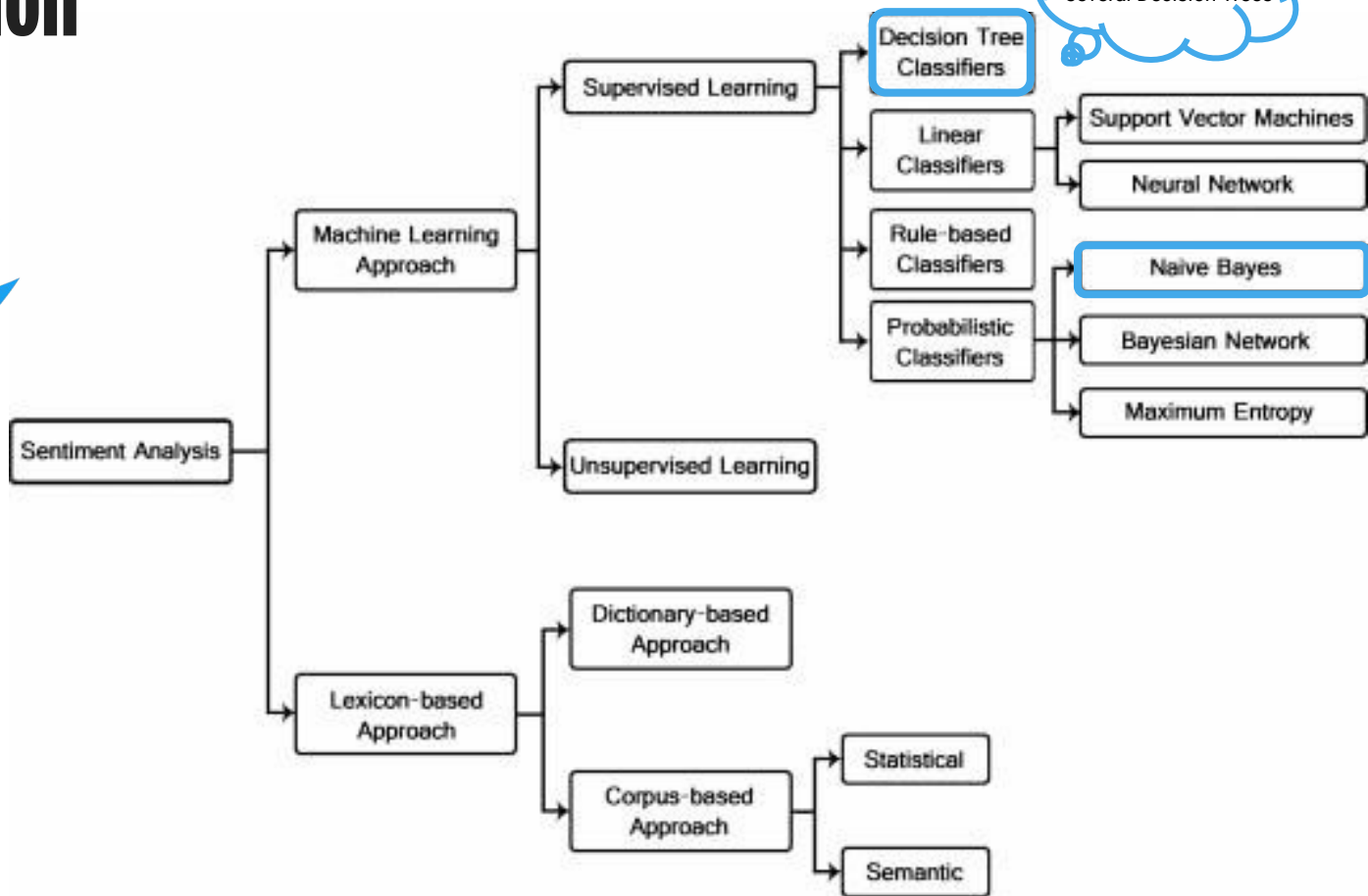
# Natural Language Processing (NLP)

Natural language processing (NLP) refers to the branch of computer science concerned with giving computers the ability to understand text and spoken words in much the same way human beings can.

- ❖ Pre-Processing
- ❖ Text Speech Processing
  - ➢ Tokenization
  - ➢ Stop Word Removal
  - ➢ Lemmatization
- ❖ Morphological Analysis
  - ➢ Bag-of-Words/ Count Vectorization
  - ➢ Term Frequency Inverse Document Frequency (TFIDF)
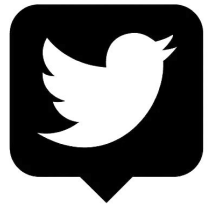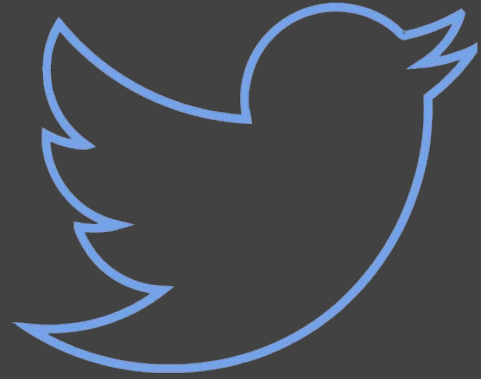  - ➢ Part-of-Speech Tagging

# Model Selection

# Model Selection

1. Random Forest: 78% accuracy

    a. Random forest is a supervised machine learning algorithm that is commonly used to solve classification and regression issues

    b. Random Forest Algorithm is that it can handle data sets with both continuous and categorical variables

2. Multinomial Naive Bayes: 80% accuracy !?!

    a. Multinomial Naive Bayes is commonly used in multinomial event models such as bag-of-words, which is a way of representing a document as vector space by counting words.

**80.6%**
**ACCURACY**

Which messages on Twitter can be **classified** as Natural Disasters, and how can we **use** that?

# How can we **use** that?

❖ Twitter is a crucial channel of communication.

❖ Crisis information shared on social media has the potential to save thousands of lives by informing others and allowing them to take preventative action.

❖ Many agencies are attempting to examine tweets programmatically to detect disasters and emergencies.

❖ This type of effort can benefit millions of people who have access to the internet and can be notified in the event of an emergency or tragedy.

# How can we use that?

❖ News organizations and disaster relief organizations are attempting to monitor tweets in real-time in order to detect calamities.

❖ Millions of ordinary people would benefit from this, as it would help them avoid potential disasters.

❖ People could take preventative action if they were notified in real-time about the occurrence of disasters in a specific place.

❖ Before the situation becomes out of control, government agencies can carry out an evacuation too.

# Questions & Comments