



NLP DISASTER TWEETS

A binary classification problem
By Cristina Sahoo

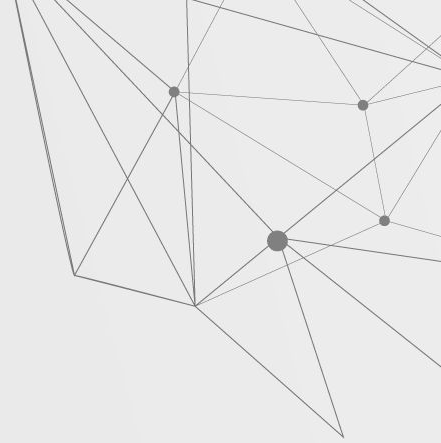
PROBLEM STATEMENT

Kaggle: In this competition, you're challenged to *build a machine learning model that predicts which Tweets are about real disasters and which one's aren't*. You'll have access to a dataset of 10,000 tweets that were hand classified.



WHY IS THIS IMPORTANT?

Twitter has become an important communication channel in times of emergency. The ubiquitousness of smartphones enables people to announce an emergency they're observing in real-time. Because of this, more agencies are interested in programatically monitoring Twitter (i.e. disaster relief organizations and news agencies).



DATA DICTIONARY

Feature	Type	Description
id	int64	a unique identifier for each tweet
text	object	the text of the tweet
location	object	the location the tweet was sent from (may be blank)
keyword	object	a particular keyword from the tweet (may be blank)
target	int64	in train.csv only, this denotes whether a tweet is about a real disaster (1) or not (0)

WORD CLOUDS

A word cloud centered around the word 'Deeds'. The word 'Deeds' is the largest and most prominent, written in a dark blue font. Other words are arranged around it in various sizes and colors, including green, yellow, and blue. The words include: shelter, wildfires, Ronge, place, got, earthquake, residents, Reason, evacuation, LaTwo, Canada, cranes, giant, receive, asked, M, near, people, Ruby, Forest, Alaska, sent, photo, and Sask.

A word cloud centered around the word 'love'. The word 'love' is the largest and most prominent, written in a dark blue font. Other words are arranged around it in various sizes and colors, including green, yellow, and blue. The words include: fast, fruits, Cramer, wrecked, lovely, Li, Disney, ready, British, Great, boxes, Summer, atmosphere, Iger, explode, goooooooooaaaaa, words, car, and engineshed.



MODELS

Models used:

- Dense Neural Network with three layers, 16 units each, activation function 'relu', and output function 'sigmoid'; optimizer 'rmsprop', loss function 'binary crossentropy', and metric 'accuracy'. Accuracy 79%.
- Multinomial Naive Bayes. Accuracy 80%.





