

# Disaster or Not ?

Tweet classification

Big-Scale analytics



**Team Google**

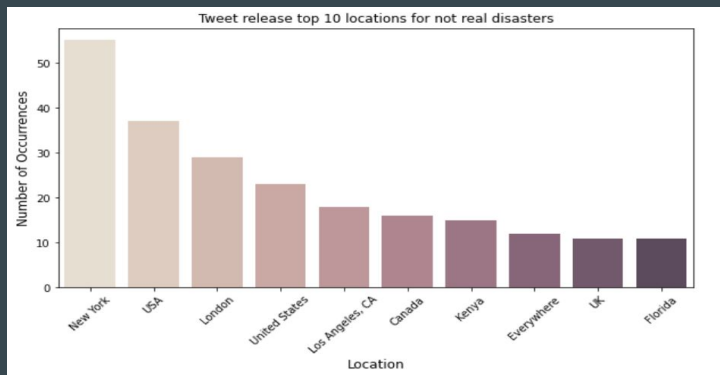
Florian Emery

Ibrahim Ounon

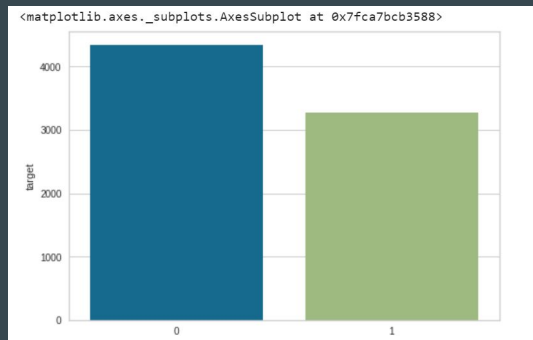
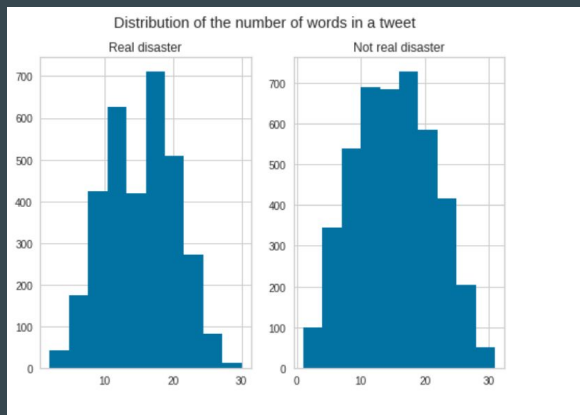
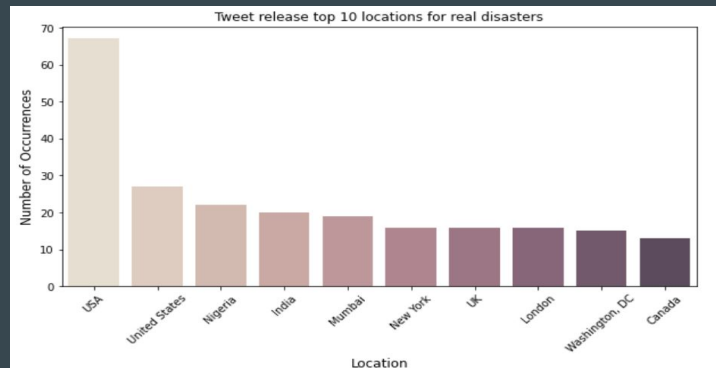
Pau Gallardo Campos

Sarah Büchner

# Data exploratory



VS



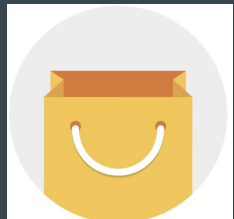
*During the data exploratory phase, we did not find anything very significant that could lead us to build on our models.*

# Bags of words - Word frequency

⇒ Creation of a basic word bag plus addition of frequently used words in the tweet to optimize our model

added a binary column "dis\_matches". If in the tweet there is a match with our list then define 1 otherwise 0

List of disaster words



List of frequent words  
in tweet about a  
disaster

Suicide
flood
Volcano
...



List of frequent words  
in tweet **not** about a  
disaster

Soul
Bag
Fruits
...



**Final** List of disaster words



# Decompressing abbreviation

[illegible]

# Data cleaning: Studying duplicates

Ambiguity of 305 duplicates -> reclassification based on avg target

	Number of records in train set	Target mean
text		
He came to a land which was engulfed in tribal war and turned it into a land of peace that is . Madinah . # ProphetMuhammad # islam	6	0.333333
The Prophet ( peace be upon him ) said 'Save yourself from Hellfire even if it is by giving half a date in charity . '	6	0.333333
To fight bioterrorism sir .	4	0.500000
.POTUS # StrategicPatience is a strategy for # Genocide ; refugees ; IDP Internally displaced people ; horror ; and so on . <a href="https://t.co/rqWuoy1fm4">https://t.co/rqWuoy1fm4</a>	4	0.750000
that horrible sinking feeling when you've been at home on your phone for a while and you realise its been on 3G this whole time	4	0.500000

```
if avg_target <= 0.5:  
    train['target'] = 0  
else  
    train['target'] = 1
```

# Data cleaning: master cleaning function

## Tokens removed

- Urls
- Html tags
- Tokens not containing alphabet characters
- Tokens with equal or less than 2 characters
- Stopwords
- Punctuation

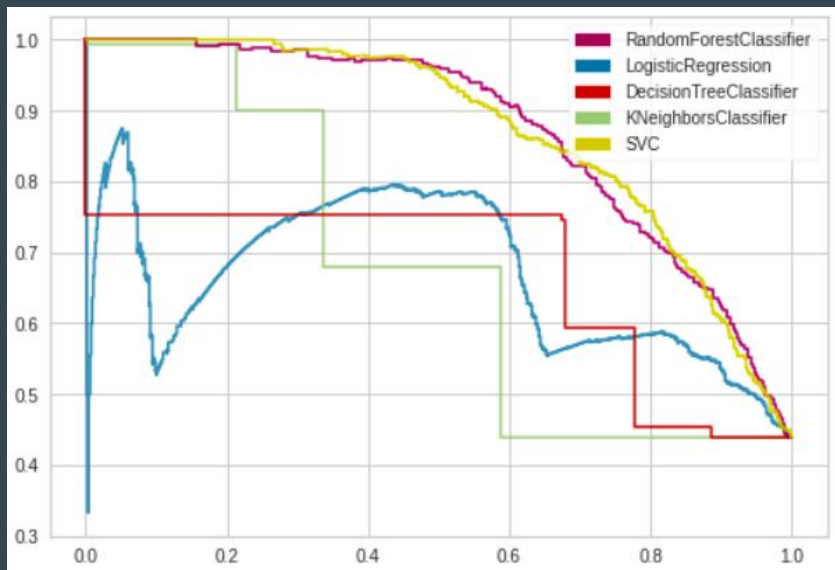
## Tokens modification

- Stemming
- Lemmatizing

## Arguments

- Array containing the tweets.

# Sklearn used methods



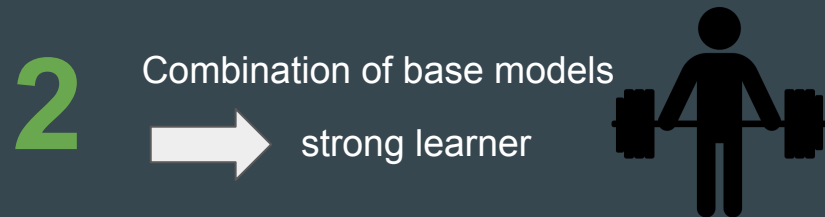
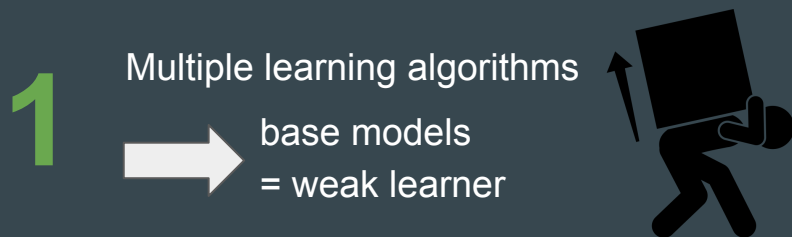
## Using Countvectorizer

Accuracy of LogisticRegression is 0.7997373604727511  
Accuracy of KNeighborsClassifier is 0.6900853578463558  
Accuracy of DecisionTreeClassifier is 0.7774130006565988  
Accuracy of RandomForestClassifier is 0.8063033486539725  
Accuracy of SVC is 0.8036769533814839

## Using TfidfVectorizer

Accuracy of LogisticRegression is 0.7918581746552856  
Accuracy of KNeighborsClassifier is 0.757715036112935  
Accuracy of DecisionTreeClassifier is 0.7321076822061721  
Accuracy of RandomForestClassifier is 0.7951411687458962  
Accuracy of SVC is 0.8108995403808273

# Ensemble Methods



Correct combination of base models  
=> more accurate and/or robust models



Improve predictive performance



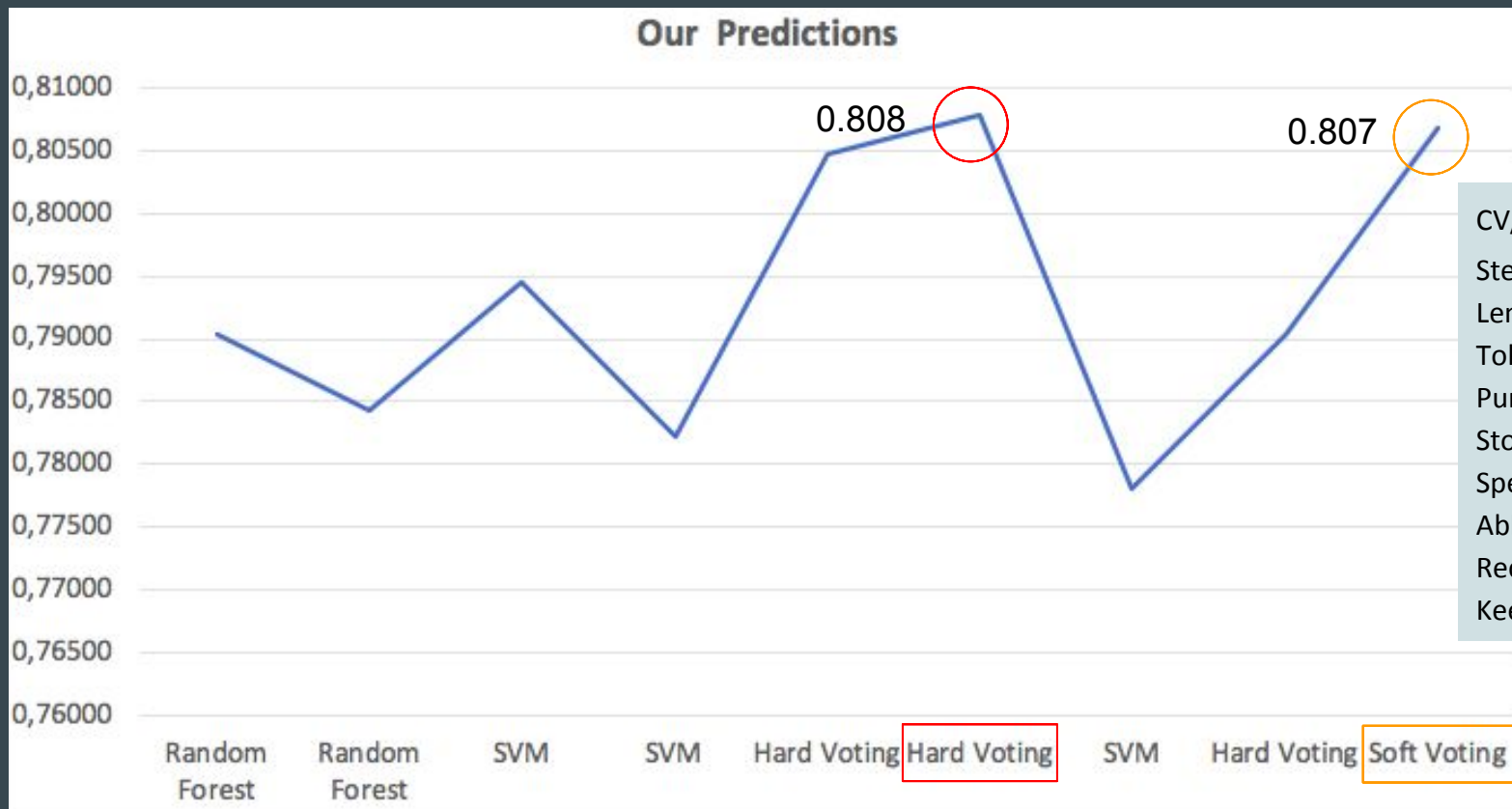
Produces more accurate solutions  
than a single model would



Higher computation time  
Reduced model interpretability  
Not always better results



# Submissions in Kaggle



CV/Tfidf Vectorizer  
Stemming  
Lemmatization  
Tokenization  
Punctuation removal  
Stopwords removal  
Special characters removal  
Abbreviations  
Reclassified wrong targets  
Keep hashtag words

# Conclusion

Ensemble methods



more accuracy



Not always!



Importance of using the correct model



Higher computation time

Reduced model interpretability



Importance of data cleaning and tokenization

**Disaster or not ?**