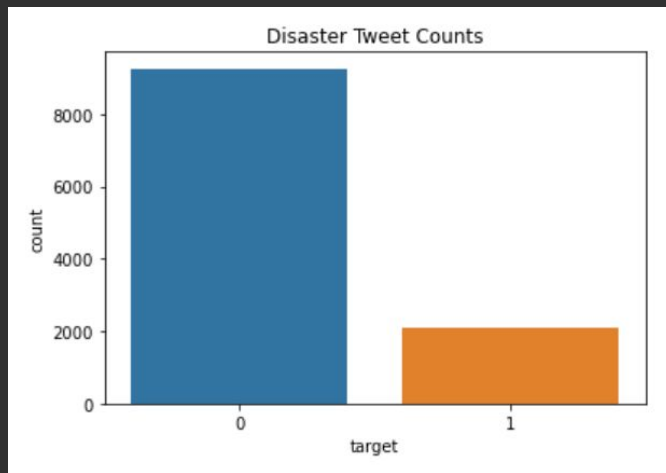# Disaster Tweet Detection

Justin Bell

How can we build an NLP model that is able to classify tweets as disasters with high accuracy?
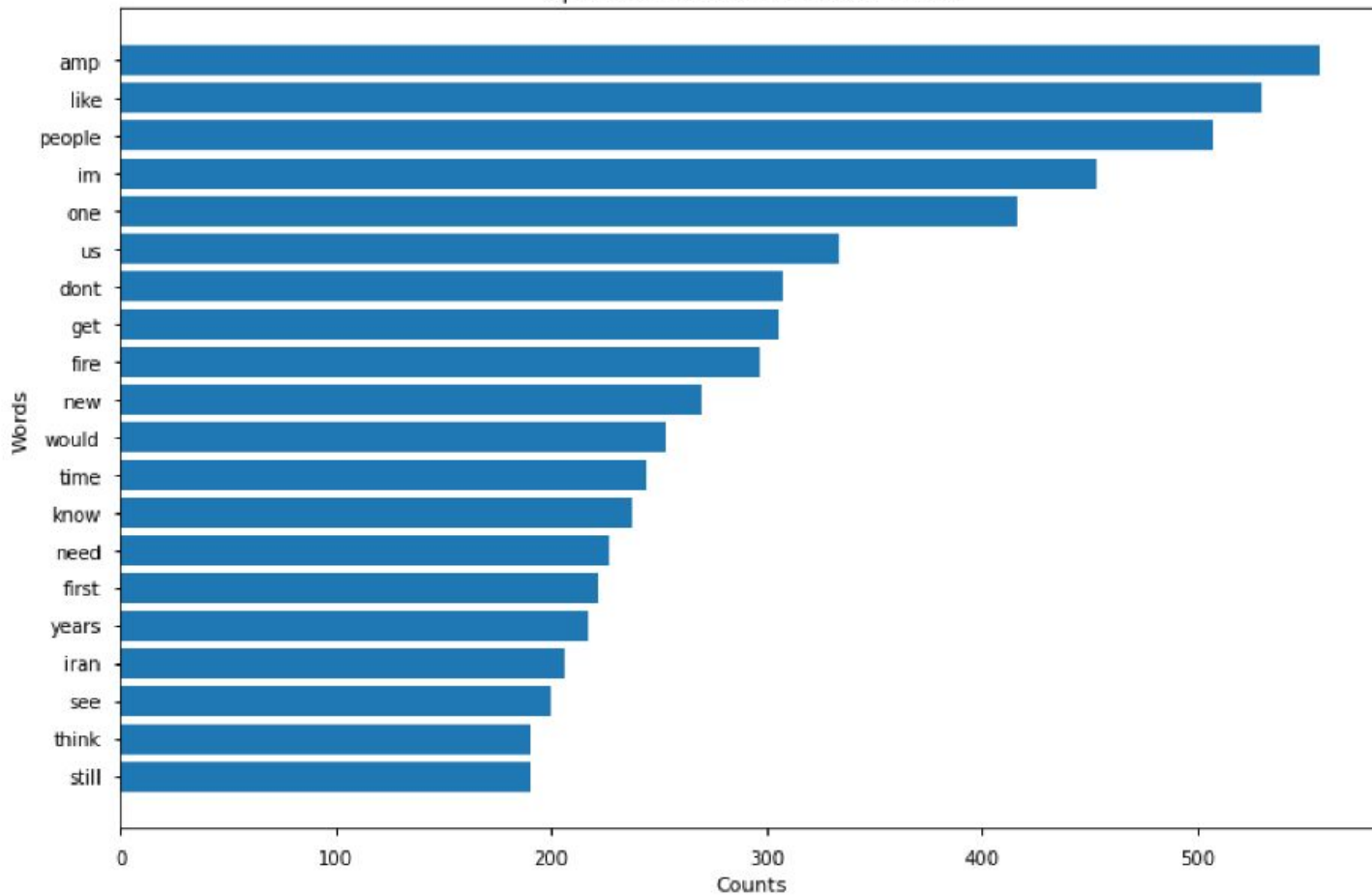
# Data Wrangling

- Import data
- Check for balanced target classes
- Data cleaning
  - Converting text to lowercase
  - Regular expression string removal
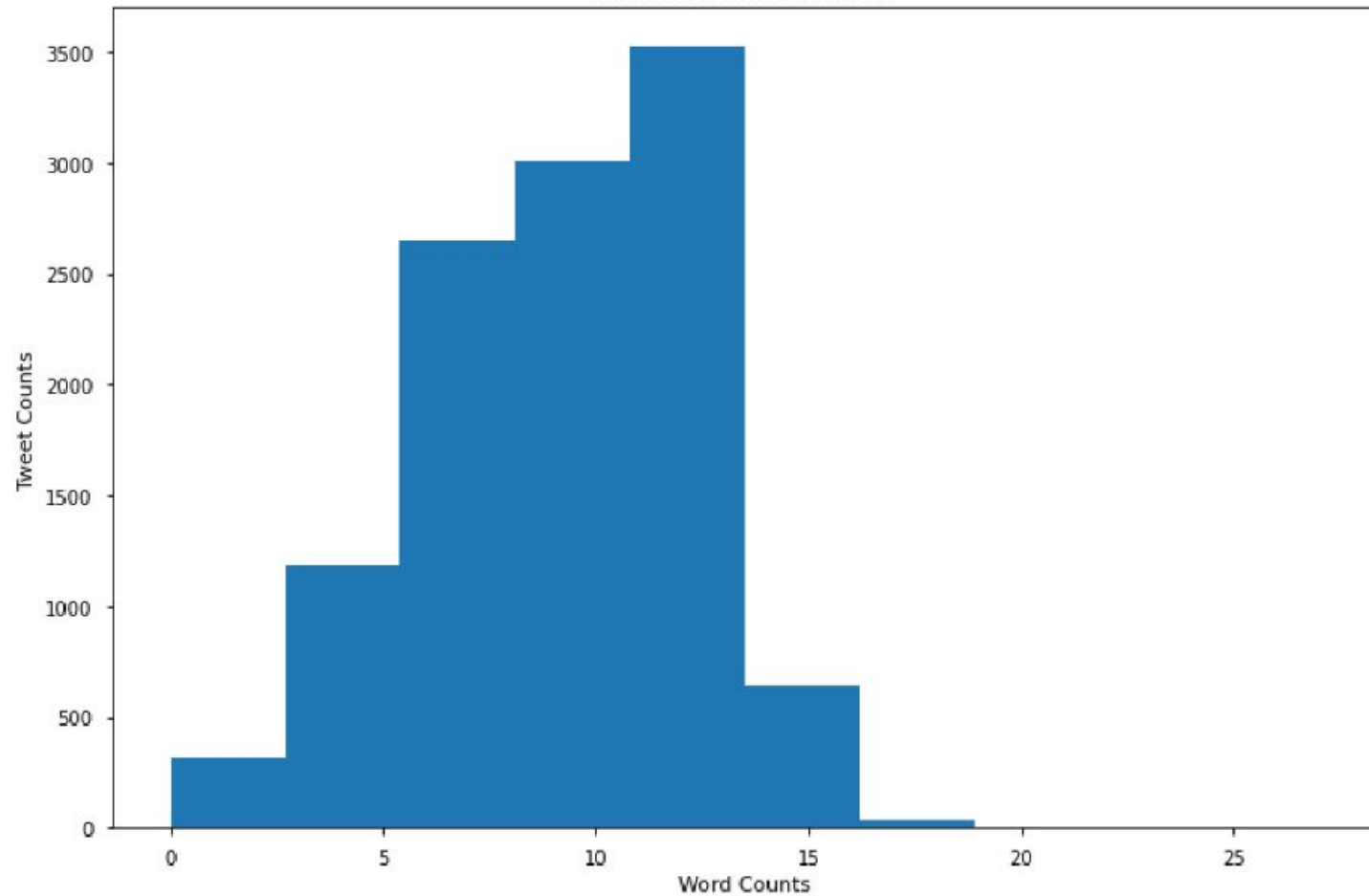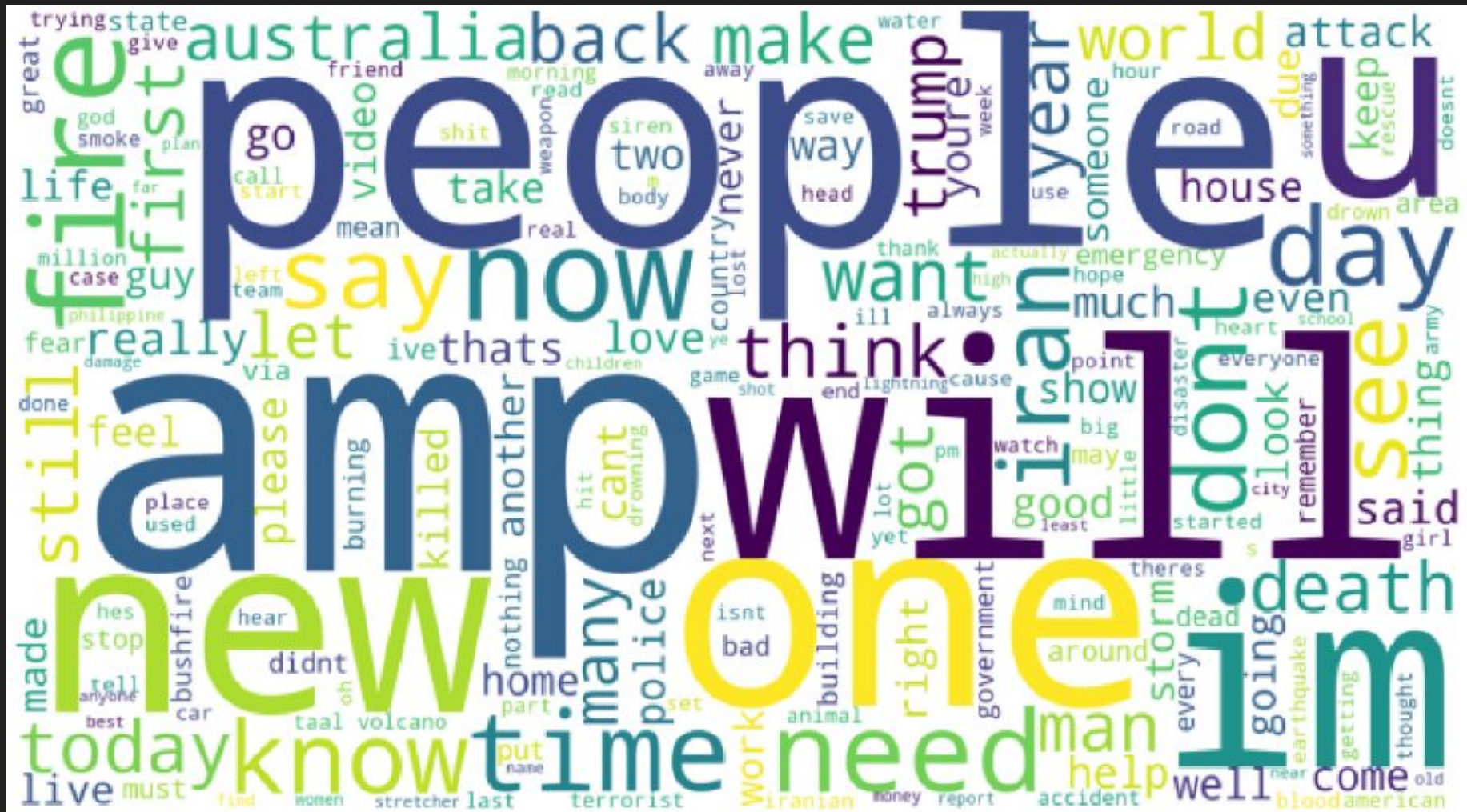  - Tokenized text
  - Remove stop words

# EDA

- Create some plots:
  - Top 20 Words → hbar plot
  - Words per Tweet → histogram
  - Word Cloud

Top 20 Common Words Across Tweets

# Preprocessing

- Remove unused columns
- Remove one null row
- Split data into train/test sets
  - 80/20
- Create vectorizers
  - Count
  - TFIDF
- Fit with training features (tweets)
- Transformed training and test features

# Modeling

- Tried 5 different classification models
    - Logistic Regression
    - Random Forest
    - KNN
    - SVM
    - Gradient Boosting

# Model Hyperparameters

## Logistic Regression:

- penalty = l2
- max_iter = 500
- C = 0.1
- solver = lbfgs

## KNN:

- n_neighbors = 6
- n_jobs = -1

## Gradient Boosting:

- n_estimators = 100
- max_depth = 10
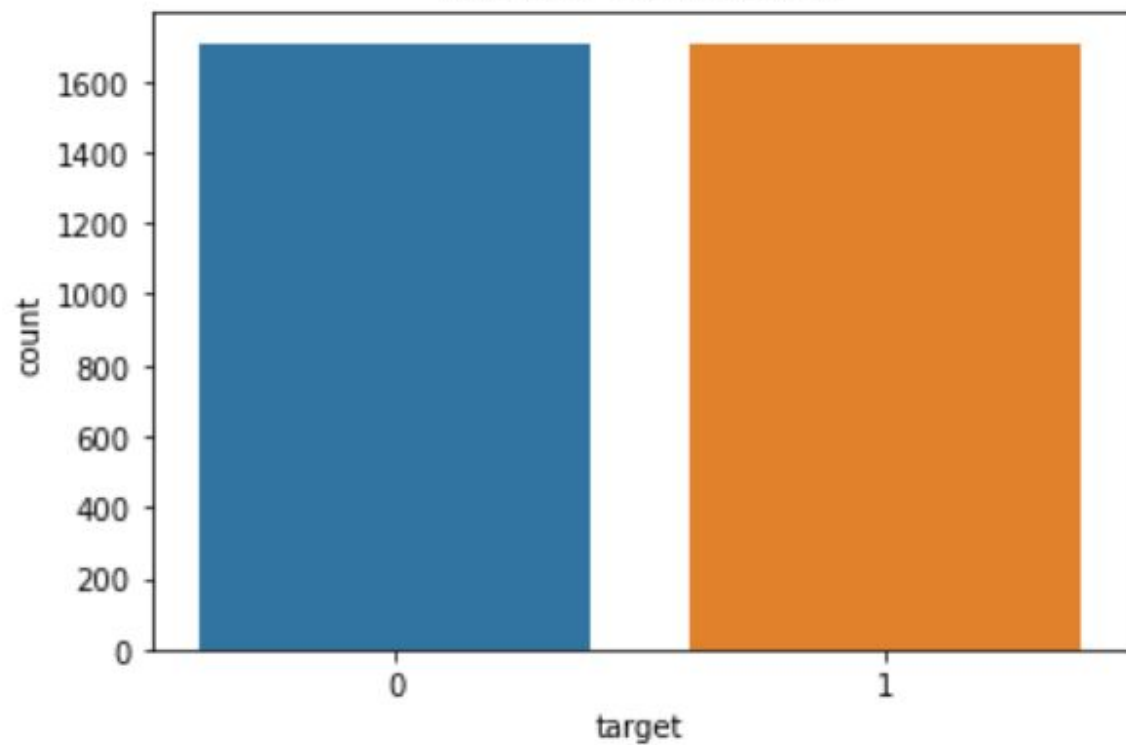- max_features = auto
- learning_rate =0.1

## Random Forest:

- n_estimators = 500
- n_jobs = -1

## SVM:

- C= 1
- gamma= 1

| Vectorizer | Model | Recall |
| --- | --- | --- |
| Count | Logistic Regression | **0.54** |
| TFIDF | Logistic Regression | 0.36 |
| Count | Random Forest | 0.49 |
| TFIDF | Random Forest | 0.49 |
| Count | KNN | 0.20 |
| TFIDF | KNN | 0.18 |
| Count | SVM | 0.19 |
| TFIDF | SVM | 0.44 |
| Count | Gradient Boosting | 0.36 |
| TFIDF | Gradient Boosting | 0.36 |

Disaster Tweet Counts

| Imbalanced/Balanced | Vectorizer | Model | Recall |
|---|---|---|---|
| Imbalanced | Count | Logistic Regression | **0.54** |
| Imbalanced | TFIDF | Logistic Regression | 0.36 |
| Imbalanced | Count | Random Forest | 0.49 |
| Imbalanced | TFIDF | Random Forest | 0.49 |
| Imbalanced | Count | KNN | 0.20 |
| Imbalanced | TFIDF | KNN | 0.18 |
| Imbalanced | Count | SVM | 0.19 |
| Imbalanced | TFIDF | SVM | 0.44 |
| Imbalanced | Count | Gradient Boosting | 0.36 |
| Imbalanced | TFIDF | Gradient Boosting | 0.36 |
| Balanced | Count | Logistic Regression | 0.75 |
| Balanced | TFIDF | Logistic Regression | 0.77 |
| Balanced | Count | Random Forest | **0.78** |
| Balanced | TFIDF | Random Forest | 0.75 |
| Balanced | Count | KNN | 0.26 |
| Balanced | TFIDF | KNN | 0.20 |
| Balanced | Count | SVM | 0.29 |
| Balanced | TFIDF | SVM | 0.76 |
| Balanced | Count | Gradient Boosting | 0.72 |
| Balanced | TFIDF | Gradient Boosting | 0.70 |

Most Important Words for Classification

# Future Scope

- Best model can be used by disaster relief organizations
- Explore more models
  - Deep Learning Neural Networks
- Plot word clouds for subsetted data (TP, TN, FP, FN)
- Similar projects