# Analysis

## Dataset analysis

The dataset contains target classes of real and fake labeled zero and one with the corresponding tweets.

Our task includes using these tweets over natural language processing to get the corresponding analysis.

## Preprocessing and Visualization

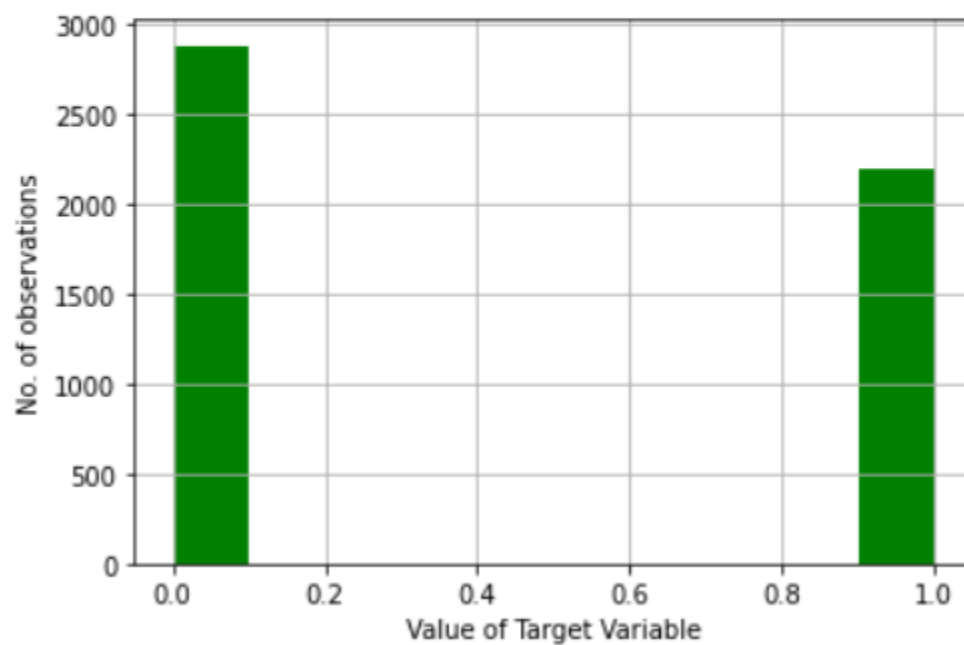Count of each target variable is being visualized through this figure.

```
Distribution of the target:
0    2884
1    2196
Name: target, dtype: int64
```
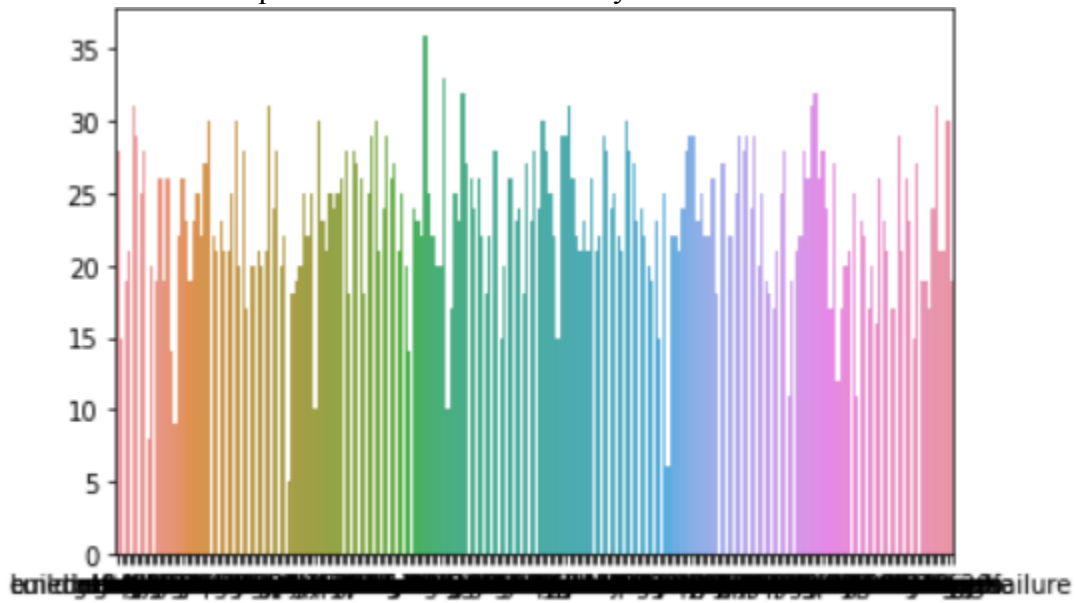
### Target distribution

After this we plotted the count of each keyword.



Steps involved in Preprocessing

1. Removing the null values.
2. Removing Double Spaces, Hyphens and arrows, Emojis, URL, another Non-English or special symbol
3. Replacing wrong spellings with correct ones.
4. Removing all columns except text and target.
5. Splitting the dataset into train and test sets.

The length of tweets for the false target is always greater than that of the real one . We can draw a conclusion that we will be getting more unique words of tweets for treal targets hence there might be a greater prior probability for the same real target.

We also analysed the the correlation between the length of tweet and the target

```
Maximum Tweet Length: 132
Minimum Tweet Length: 0
Average Tweet Length: 37.70393700787402


Correlation 16.166298308946054
```
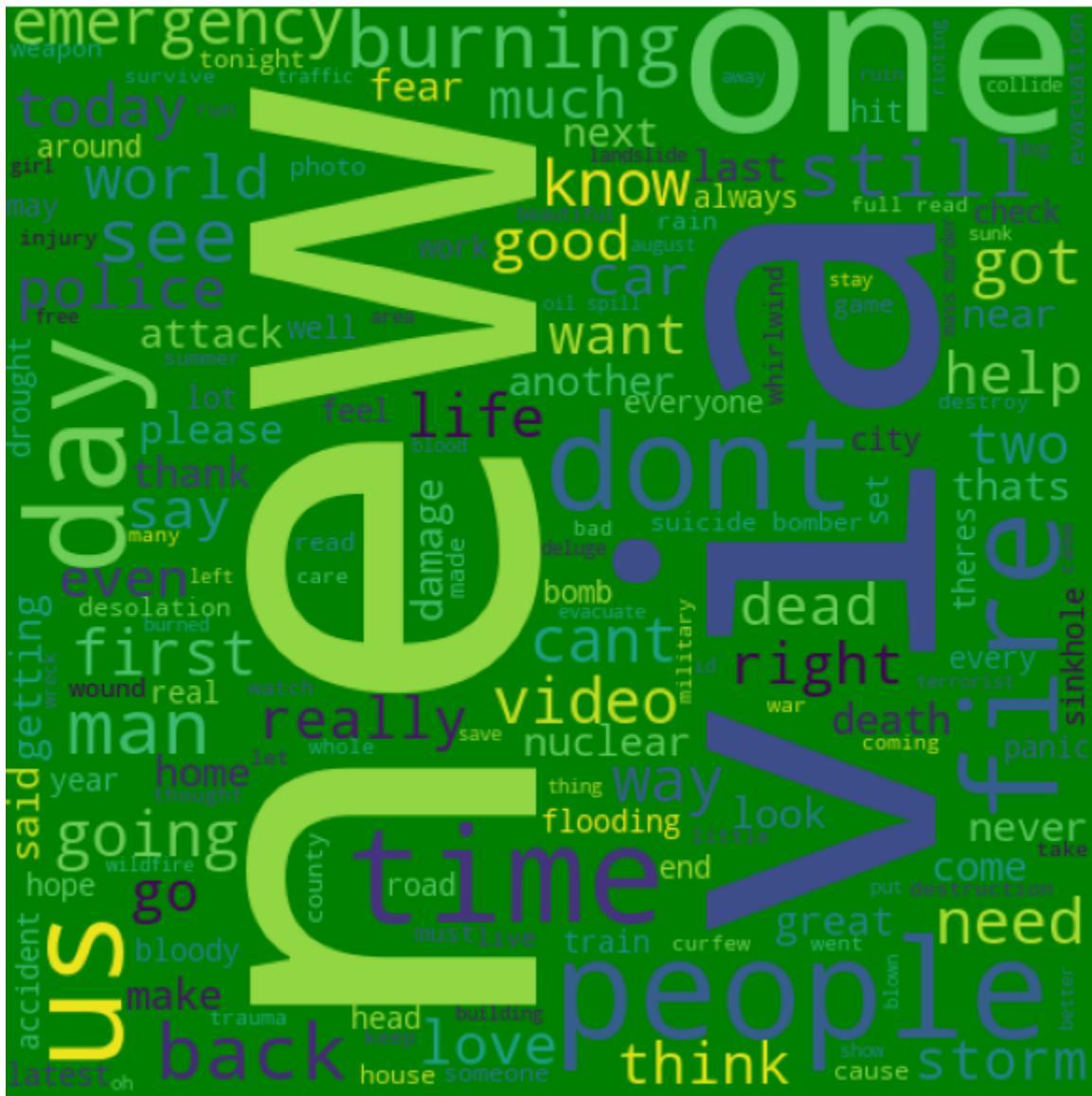
## Word Cloud

A word cloud is a collection, or cluster, of words depicted in different sizes. The bigger and bolder the word appears, the more often it's mentioned within a given text and the more important it is.

Target 0

Target 1

Does the sum of the unique words in target 0 and 1 sum to the total number of unique words in the whole document? Why or why not? Explain in the report.

```
Number of unique words in class 0 : 4194
Number of unique words in class 1 : 3206
sum =  7400
Total unique words = 5380
```

As you can see, the number of unique words in target 0 + number of unique words in target 1 is more than the total number of unique words in the whole document. It is because some words are common for both targets.

## Applying Bayes

$$P(positive|review) \ \alpha \ P(w_1|positive) * P(w_2|positive) * P(w_3|positive) * P(w'|positive) * P(positive)$$

We have calculated the probability of occurrence of the word in a class, we can now substitute the values in the Bayes equation.

After doing the predictions we calculated the True Positive, True Negative, False Positive and False Negative for the dataset which we then used to calculate the precision, recall, accuracy, etc.