

Capstone 2 Milestone Report 1

Toxic Comment Challenge

By Minh Ngoc Pham

Data Science Career Track

Disclaimer: the dataset and this project for this competition contains text that may be considered profane, vulgar, or offensive.

Problem Statement

There is an increasingly high number of toxic behaviour and comments over the internet which are making it difficult to have meaningful discussion on the net. As the world is becoming more and more technologically involved and people demonstrate more presence on the internet compared to in real life, it is important to make it a healthy and safe place to express opinions.

One does not need to look far but at Youtube, where presence of toxic and hatred comments is ubiquitous in almost any video. This has led many Youtube channels to disable the comment section, which in itself might be a bit of an extreme measure given the fact that genuine commentators might be interested in the content and want to engage in meaningful discussion. The need for a mechanism to deal with outright toxic and unhealthy comments are more than ever relevant in this day and age.

The Dataset

Dataset: <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/>

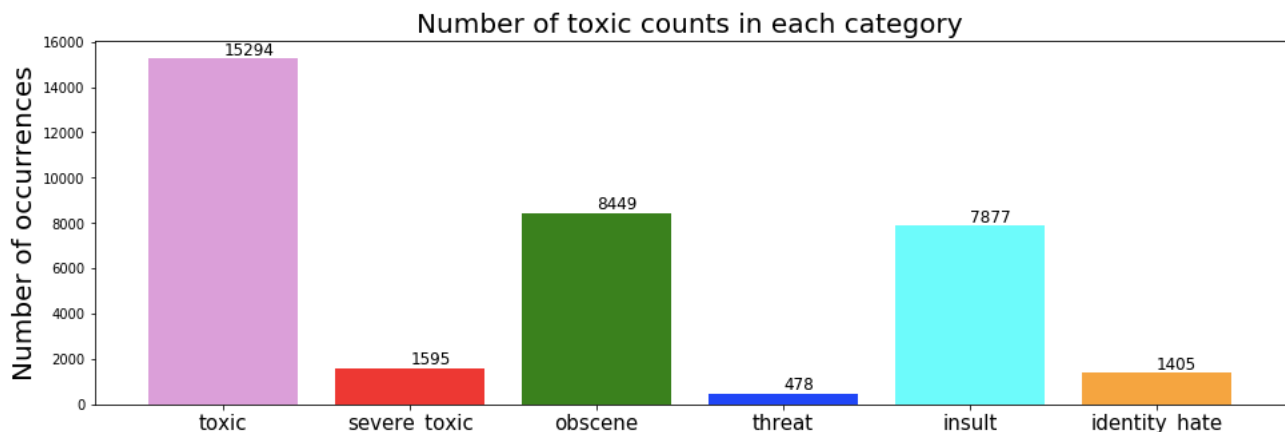


The dataset chosen for addressing this issue is the Toxic Comment Classification Challenge by Kaggle that was published in March 2018. Conversation AI is a research initiative by Jigsaw and Google aiming at improving online conversation. They have created a lot of publicly available models but there are still errors involved. This competition challenged others to find a model that can better detect negative online behaviour.

This dataset contains Wikipedia comments that have been rated for toxic behaviour. There are 6 classes for 6 types of toxic behaviour including: toxic, severe_toxic, obscene, threat, insult, identity_hate. There are 159,571 comments in the training set (which also includes comments that are rated as neutral, meaning being scored 0 for all classes). The test set includes 153,164 comments.

Initial Data Visualisation

We will first have a look at our dataset to see how many comments being classified as a certain label:



We can see from the above graph that “toxic” category has the highest number of comments with 15,294 comments classified as toxic. The “threat” class has the lowest number of comments with only 478 comments in this class. We notice that the whole training dataset contains 159,571 comments in total. With this in mind, the class with the highest number of comments (“toxic”) account for even less than 10% of all the comments. This suggests a class imbalance issue with a very high number of clean comment.

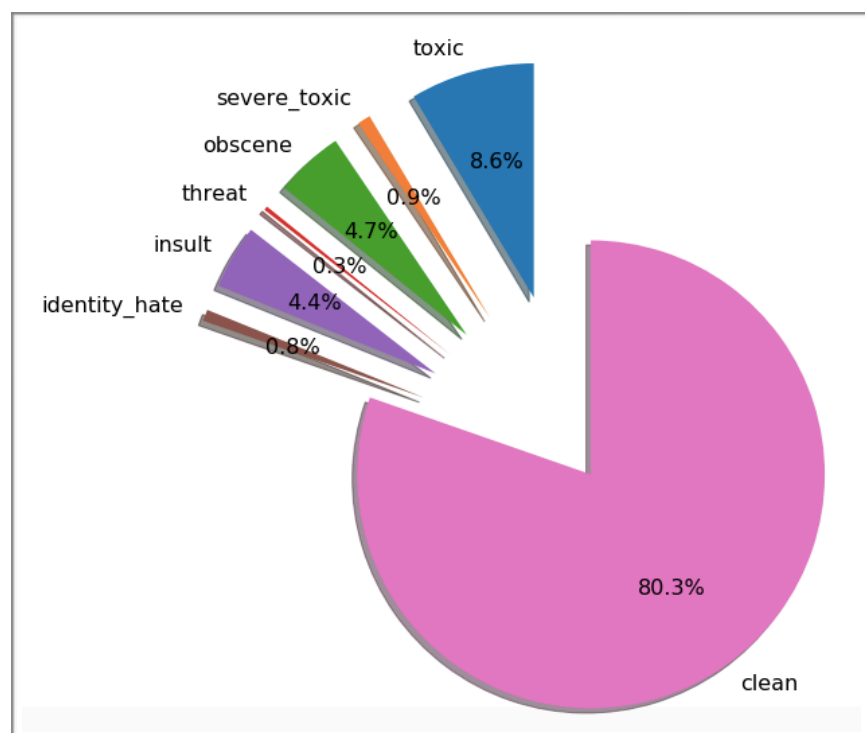
To further shed light on this issue, we will create a new column classifying a comment as clean or not.

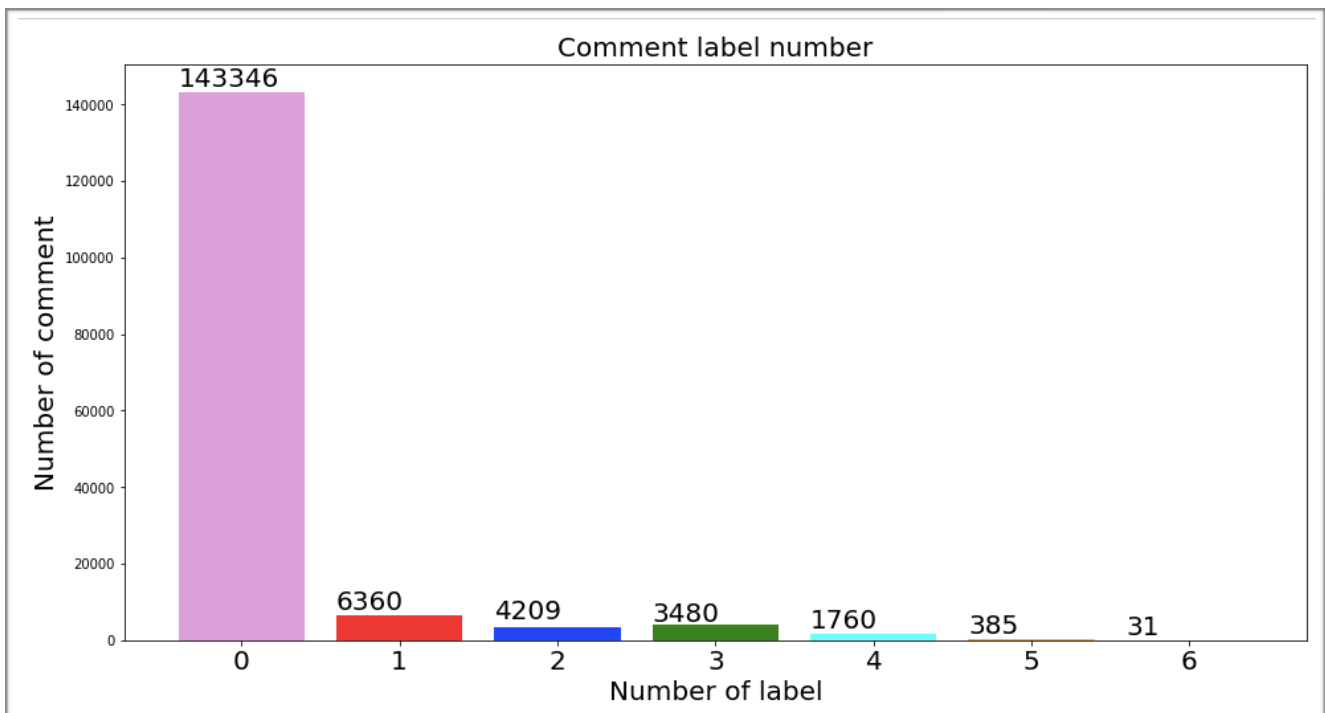
The pie chart shows that we have about 80.3% of the comments that do not belong to any of the 6 classes.

We might face a class imbalance situation in this case. As a result, accuracy score might not be the most meaningful determiner of how good a model is. We might need to consider: ***confusion matrix, precision, f1 score and recall.***

Checking the number of labels per comment

Each comment can be classified as more than one label. They can have many label classification simultaneously with the highest number of label of 6 and the lowest of 0.





The bar chart shows the count of comments with certain number of labels:

As expected, we see a high number of comment with 0 label ('clean' comment) of 143,346 comments. It is then followed by the count of comment with 1 label (6,360 comments). There are only 31 comments that are classified as all 6 of the labels.

Data Cleaning and Wrangling

We will have a look at the train data and the test data to see if there is any null value:

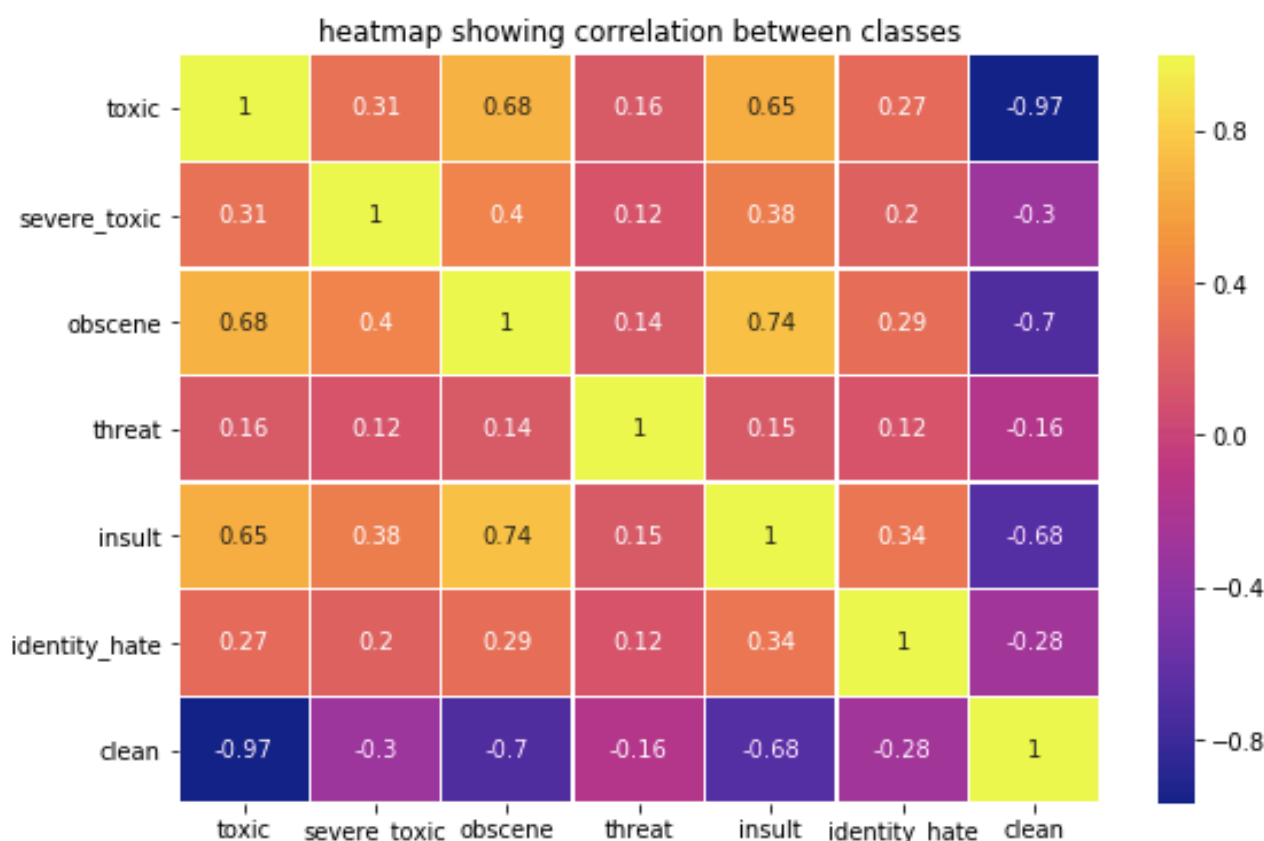
As expected, there is no null values in any of the set. As this is a Kaggle competition dataset. The data was probably already clean.

- ◆ There is no null value in either the Train set or the Test set

```
Checking for null values in the Train set
comment_text    0
toxic            0
severe_toxic    0
obscene         0
threat          0
insult          0
identity_hate   0
dtype: int64
Checking for null value in the Test set
comment_text    0
dtype: int64
```

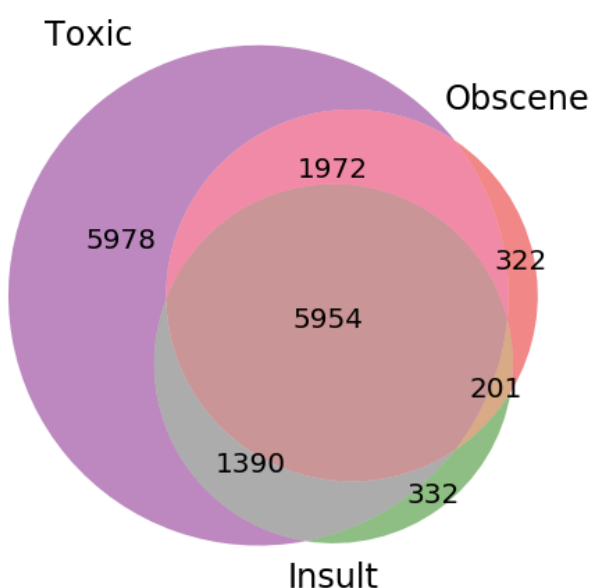
Testing for correlation between feature variables

As we can see, the 7 feature variables (including clean) are dichotomous categorical variables (with values as either 0 or 1) so we can apply the Pearson correlation (also known as point-biserial correlation coefficient) in this case. We will do a correlation heatmap in order to look at the correlation between these variables.



Code notice: the .corr() formula simply ignore the comment_text variable as it is a non-numerical variable

- ◆ There is high correlation between clean and toxic (negatively correlated) suggesting that the comments are likely to be classified as either toxic or clean



- ◆ obscene and insult categories also have high correlation with toxic (0.68 and 0.65 respectively). Perhaps we can explore further with a Venn diagram to see

It is interesting to note here that:

- ◆ The majority of comments classified as obscene are also classified as toxic (1972 + 5954 = 7926 comments out of 8449 obscene comments which is roughly 93.8%)
- ◆ and the majority of comments classified as insult are also classified as toxic

(5954+1390=7344 comments out of the 7877 insult comments which is roughly 93.2%)

- ✦ 5954 comments are classified as obscene, insult and toxic

Text Cleaning and Word Clouds

An initial look at the `comment_text` columns show that there are some characters in the text that can be cleaned and preprocessed. The following steps were taken to clean the texts:

- ◆ Combining the train data and the test data for cleaning (using `pd.concat()`)
- ◆ Removing all punctuation
- ◆ Removing the `\n` in the comments
- ◆ removing the digits in the comments
- ◆ Splitting combined words
- ◆ Converting words into lowercase
- ◆ tokenize comments
- ◆ Setting and removing stopwords (`set(stopwords.words('english'))`)
- ◆ Converting words in base form/lemmatize
- ◆ Split the data back in clean training data and clean test data

Once the text is nice, clean and tokenise. We can perform Word Cloud to have a look the words that appear a lot in certain labels.



