

Toxic Comment Classification

Kuan - Han Chen
July 15th, 2018

Definition

Project Overview

Discussing things you care about can be difficult. The threat of abuse and harassment online means that many people stop expressing themselves and give up on seeking different opinions. Platforms struggle to effectively facilitate conversations, leading many communities to limit or completely shut down user comments.

The Conversation AI team, a research initiative founded by Jigsaw and Google (both a part of Alphabet) are working on tools to help improve online conversation. One area of focus is the study of negative online behaviors, like toxic comments (i.e. comments that are rude, disrespectful or otherwise likely to make someone leave a discussion). So far they've built a range of publicly available models served through the Perspective API, including toxicity. But the current models still make errors, and they don't allow users to select which types of toxicity they're interested in finding (e.g. some platforms may be fine with profanity, but not with other types of toxic content).

In this project, I will build a multi-headed model that's capable of detecting different types of toxicity like threats, obscenity, insults. I will use naive bayes as my benchmark model. And get a baseline score through naive bayes model for my dataset, and compare with Text-CNN, which algorithm has higher AUC, ROC and accuracy.

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

Problem Statement

The Conversation AI team, a research initiative founded by Jigsaw and Google (both a part of Alphabet) are working on tools to help improve online conversation. One area of focus is the study of negative online behaviors, like toxic comments (i.e. comments that are rude, disrespectful or otherwise likely to make someone leave a discussion). So far they've built a range of publicly available models served through the Perspective API, including toxicity. But the current models still make errors, and they don't allow users to select which types of toxicity they're interested in finding (e.g. some platforms may be fine with profanity, but not with other types of toxic content).

Short-text classification is an important task in natural language processing, including selective dissemination of information to information consumers, spam filtering, or

sentiment analysis. Many different approaches have been developed for short-text classification. In machine learning models, such as Naive Bayes Classifier, Support Vector Machine(SVM), Decision Tree Classifier. In deep Neural Networks, such as Convolutional Neural Network (CNN), Long Short Term Model (LSTM), Gated Recurrent Unit (GRU), Recurrent Convolutional Neural Network (RCNN). Among these different models, Convolutional Neural Network (CNN) architecture have also demonstrated profound efficiency in NLP tasks including sentiment classification. The method of classifier we will use in this project is TextCNN, which was proposed by Y Kim in "Convolutional Neural Networks for Sentence Classification", 2014.

I will build a multi-headed model that's capable of detecting different types of toxicity like threats, obscenity, insults, and identity-based hate better than Perspective's current models. You'll be using a dataset of comments from Wikipedia's talk page edits.

Improvements to the current model will hopefully help online discussion become more productive and respectful.

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

Metrics

- ROC: The AUC value is equivalent to the probability that a randomly chosen positive example is ranked higher than a randomly chosen negative example.
- AUC: In signal detection theory, a receiver operating characteristic (ROC), or simply ROC curve, is a graphical plot which illustrates the performance of a binary classifier system as its discrimination threshold is varied.
- Accuracy: How many number of correct prediction in total number of correct predictions.
- Calculation formula :

$$\text{ROC/AUC} : \text{TPR} = \frac{TP}{TP + FN} , \text{FPR} = \frac{FP}{FP + TN}$$

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

Despite the processing by resampling, the data is still imbalanced, so I will use ROC/AUC score to measure the performance of the model.

Datasets and Inputs

The datasets are provided by Jigsaw and Google for competeion in Kaggle. Them contain training set and testing set.

- train.csv(160k rows x 8 columns) - the training set, contains 8 columns, namely id, comment_text, toxic, severe_toxic, obscene, threat, insult and identity_hate, and last six columns comments with their binary labels. train.csv visualized by pandas, shown in Fig. 1.
- test.csv(153k rows x 2 columns) - the test set, contains 2 columns, namely id and comment_text. You must predict the toxicity probabilities for these comments. To deter hand labeling, the test set contains some comments which are not included in scoring.
- test_labels.csv(153k rows x 7 columns) - labels for the test data; value of -1 indicates it was not used for scoring. It contains 7 columns, namely id, toxic, severe_toxic, obscene, threat, insult and identity_hate.

	id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
0	0000997932d777bf	Explanation\nWhy the edits made under my usern...	0	0	0	0	0	0
1	000103f0d9cfb60f	D'aww! He matches this background colour I'm s...	0	0	0	0	0	0
2	000113f07ec002fd	Hey man, I'm really not trying to edit war. It...	0	0	0	0	0	0
3	0001b41b1c6bb37e	"\nMore\nI can't make any real suggestions on ...	0	0	0	0	0	0
4	0001d958c54c6e35	You, sir, are my hero. Any chance you remember...	0	0	0	0	0	0
5	00025465d4725e87	"\n\nCongratulations from me as well, use the ...	0	0	0	0	0	0

Fig. 1. train.csv visualized by pandas

Analysis

Data Exploration and Visualization

These datasets contain a large number of text comments and classified into which type of toxicity like threats, obscenity, insults. When we import train sets and Visualize this data. Fig 1. Shows the number of each class. Then we found that data is an imbalanced data, the number of each type varies greatly. Use imbalance datasets to train your model, the accuracy measures may tell you that you have good accuracy, but the accuracy is reflecting the underlying each class exist non-uniform distributed. If we get the imbalance data, we can take following approach.

1. Change performance metric. Accuracy is not the appropriate metric to use in an imbalanced data, that would mislead us. We can use Precision, Recall, F1 – score and AUC/ROC as performance metric
2. Resample dataset. We can add copies of instances from the under-represented class, called over-sampling. And delete instances from the over-represented, called under-sampling.
3. Use different algorithm, like decision tress and SVM. These algorithm have less affected from imbalance data.

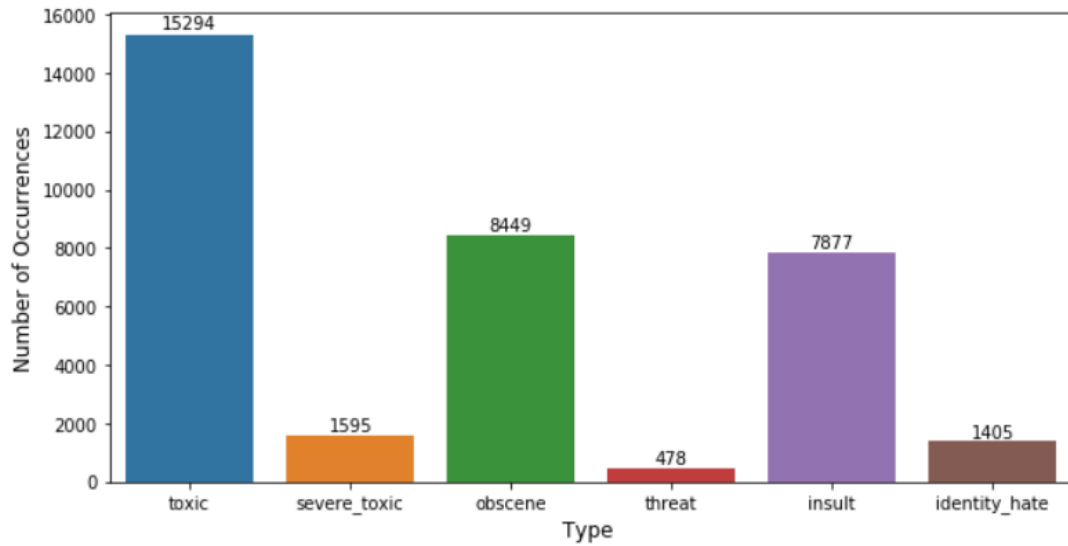


Fig 1. The number of each class

We take over-sampling to let the few materials to make it equivalent to the majority of the data. After sampling, shows in Fig 2.

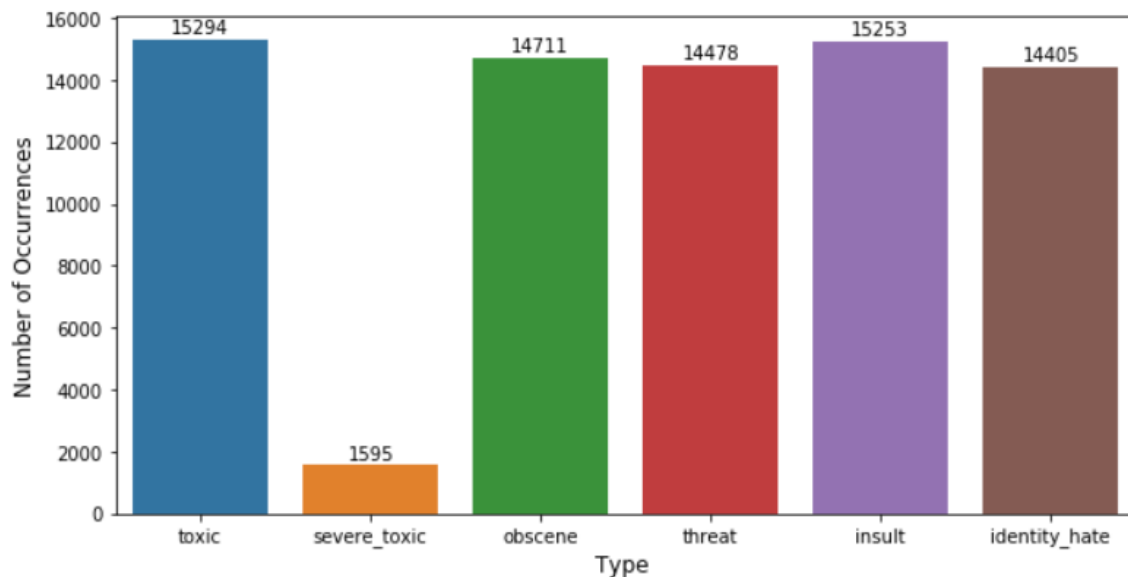


Fig. 2. After sampling, The number of each class

From Fig. 2. We still can found that the class of severe_toxic is still under-presented. We some simple way to analyze the relationship between the class of toxic and severe_toxic, shows in Fig. 3. We figure out that sever_toxic comments are always toxic, so that We increase the number of toxic and also increase the number of severe_toxic. So we don't take any action on the class of severe_toxic.

```

: print("The number of toxic get 1:",train[(train.toxic==1)].shape[0])
print("The number of severe_toxic get 1:",train[(train.severe_toxic==1)].shape[0])
print("The number of toxic and severe_toxic both get 1: ",
      train[(train.severe_toxic==1)&(train.toxic==1)].shape[0])

```

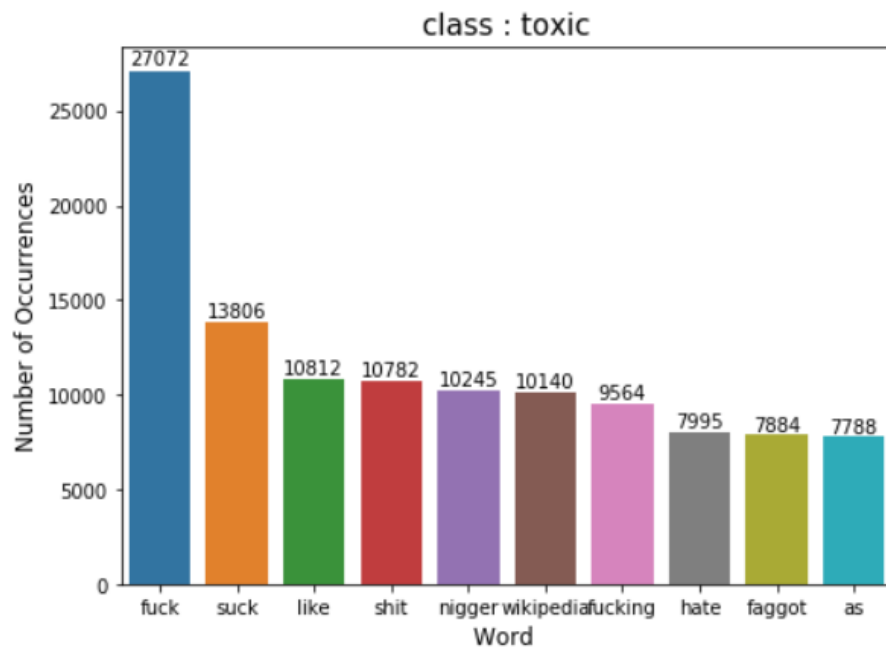
The number of toxic get 1: 15294

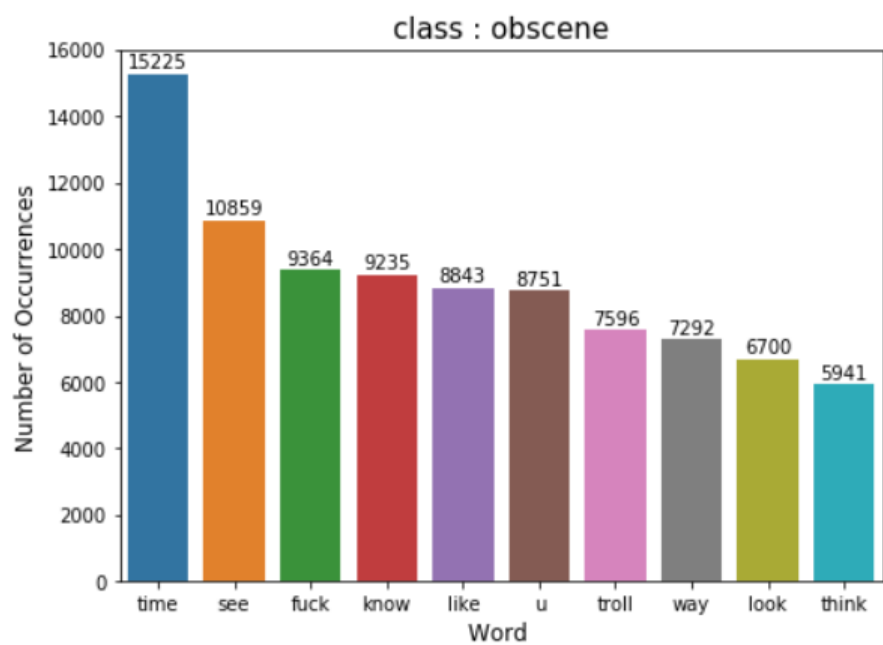
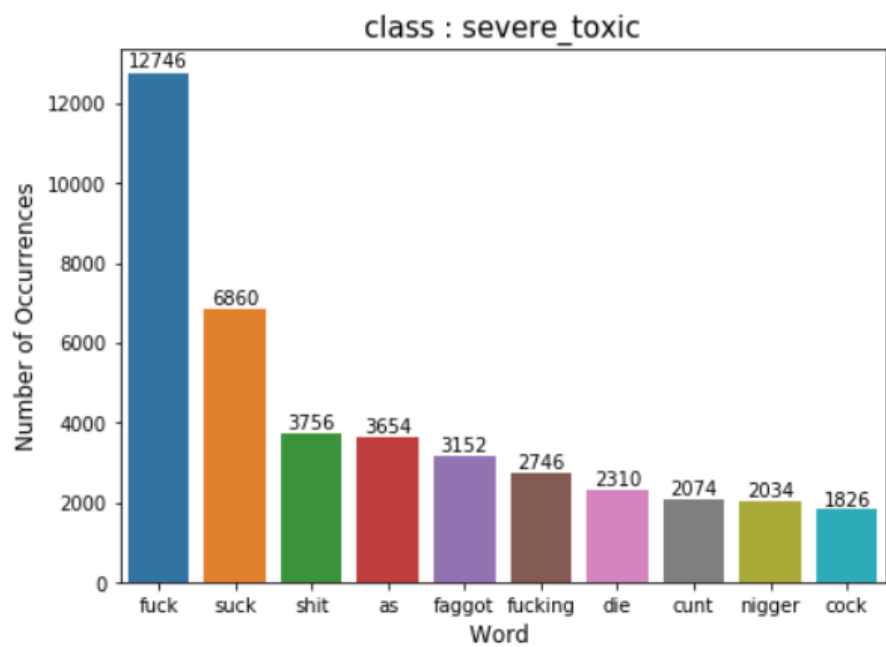
The number of severe_toxic get 1: 1595

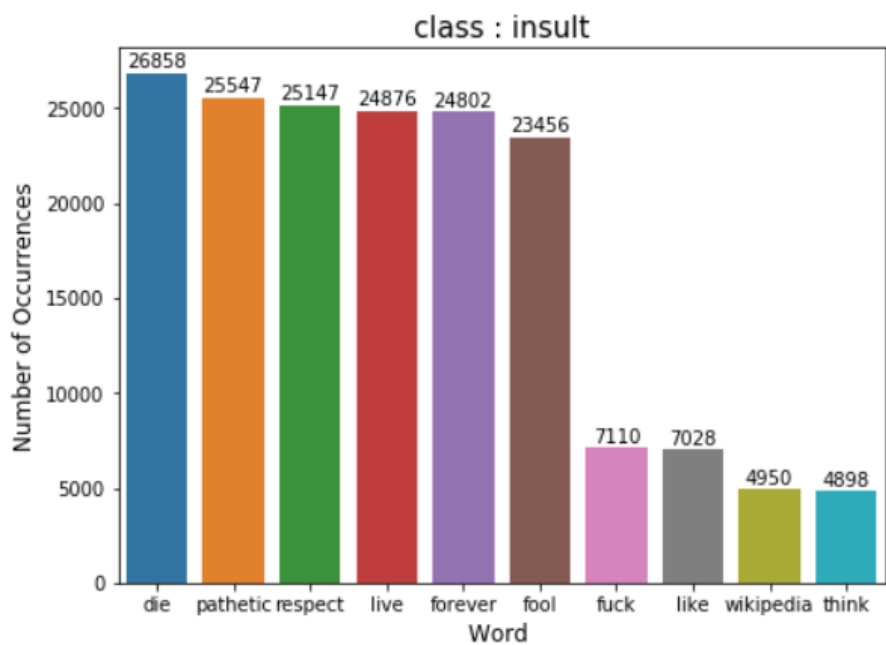
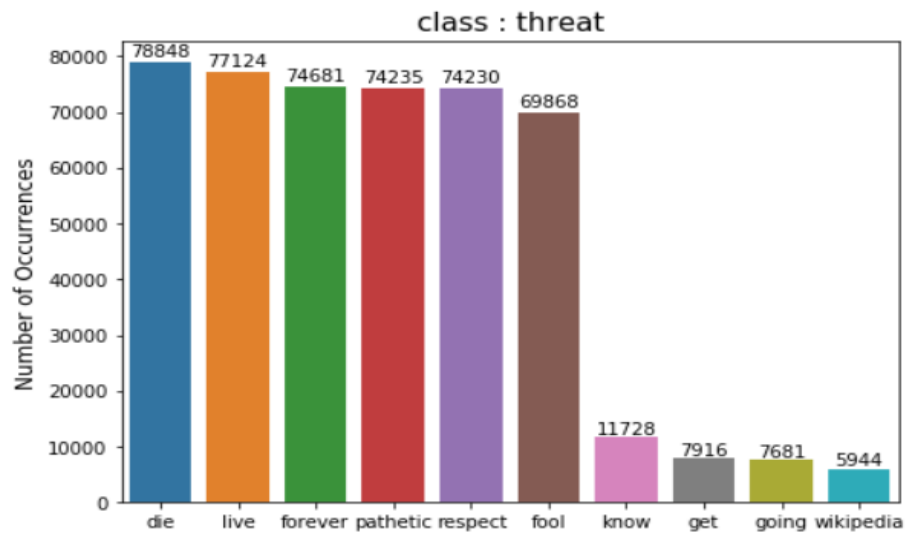
The number of toxic and severe_toxic both get 1: 1595

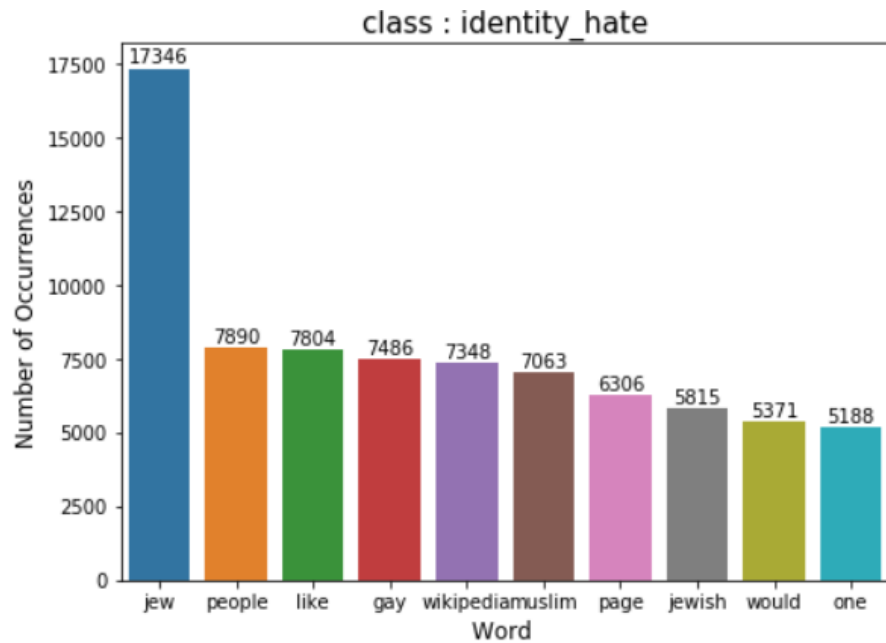
Fig. 3. analyze the relationship between the class of toxic and severe_toxic

We find out which words are often found in specific classes and visualize these top 10 words from each class are shown below.









Algorithms and Techniques

The classifier is a TextCNN, which is state in "Convolutional Neural Networks for Sentence Classification from Yoon Kim in 2014. The model architecture, shown in Fig. 4.

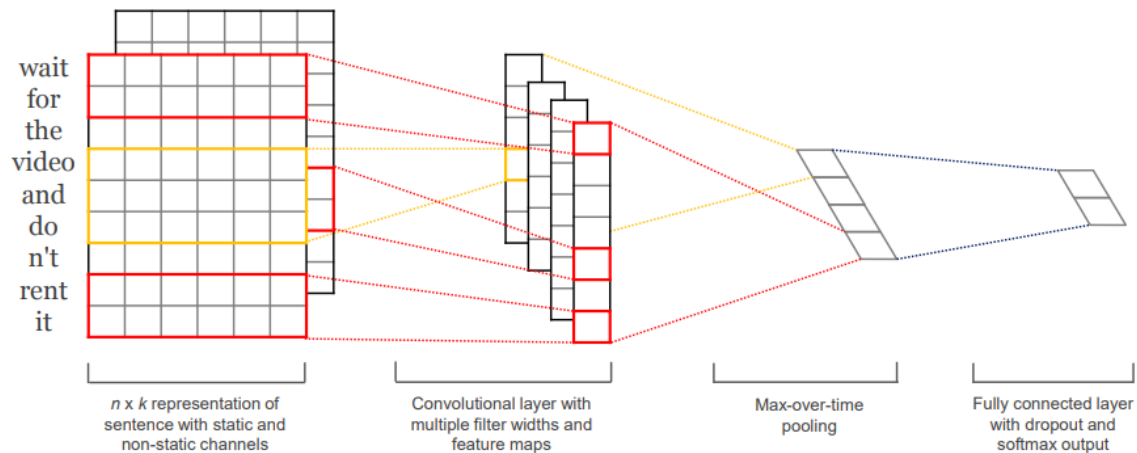


Fig. 4. Model architecture with two channels for an example sentence

Model architecture described as follows:

1. Input layer

As shown in the figure, the input layer is a matrix in which the word vectors corresponding to the words in the sentence are arranged in order (from top to bottom). Assuming that the sentence has n words and the dimension of the vector is k , then the matrix is $n \times k$.

2. Convolutional layer

The input layer obtains several Feature Maps by convolution operation. The size of the convolution window is $h \times k$, where h is the number of vertical words and k is the

dimension of the word vector. Through such a large convolution, we will get several Feature Maps with a column number of 1.

3. Pooling layer

In the pooling layer, a method called Max-over-time Pooling is used in the text. This method simply proposes the largest value from the previous one-dimensional Feature Map, which explains that the maximum value represents the most important signal. It can be seen that this kind of filtering can solve the variable length sentence, which we input. Because, no matter how many data from the Feature Map, we only need to extract the maximum value)

4. Fully connected layer + softmax

The output of the one-dimensional vector of the pooling layer is connected to a Softmax layer by means of full connection, and the Softmax layer can be set according to the needs of the task.

When we build TextCNN model, we will set some parameters. List some important parameter are as follows.

- max_features: the number of unique words.
- Maxlen: Unify the dimensions of all sentences. Make the shorter sentences has the same size with others by filling the shortfall by zeros.
- Batch size: the number of training examples in one forward/backward pass.
- Epochs: One Epoch is when an entire dataset is passed forward and backward through the neural network only once.

I will use Text-CNN as the main solution algorithm. Text-CNN has a good performance in natural language processing. The Text-CNN algorithm was published in this paper, Convolutional Neural Networks for Sentence Classification

Benchmark

To create an initial benchmark for the classifier, I use Decision Tree Classifier as benchmark model. We did not adjust any parameters from this benchmark model, we get the best accuracy is 0.846 and auc_roc score is 0.85

Methodology

Data Preprocessing

- 1. Load Data**
- 2. Data visualization**
- 3. Data Exploration and Visualization**
- 4. Resampling data(over-sampling)**
- 5. Preprocessing**

- Convert all letters to lowercase
- Apostrophe replacement
- Remove punctuation
- Remove stopwords
- Split into words
- Stem words

Implementation

1. Load Data
2. Data visualization
3. Data Exploration and Visualization
4. Resampling data(over-sampling)

I filter out the data that the class of identity_hate is equal to 1 and the class of toxic is equal to 0, random copy and add to my train set.

Here is the process of resample data. No fixed practice, you can try other combinations.

```
train1=train_app[(train_app.identity_hate==1) & (train_app.toxic==0)].sample(n=11000,replace=True)
train_app=train_app.append(train1, ignore_index=True)

train1=train_app[(train_app.threat==1) & (train_app.toxic==0)].sample(n=14000,replace=True)
train_app=train_app.append(train1, ignore_index=True)

train1=train_app[(train_app.obscene==1) & (train_app.identity_hate==1)& (train_app.toxic==0)& (train_app.insult==0)]
.train_app.sample(n=2000,replace=True)
train_app=train_app.append(train1, ignore_index=True)
```

5. Preprocessing

- Convert all letters to lowercase
- Apostrophe replacement
you're --> you are
what's --> what is
- Remove punctuation
- Remove stopwords
 - Stopwords are very common words in a language, like “the”, “is”, “she”. The value of stopwords in txt-classification is low, we will be drop them from the data. We get a set of English stop words from nltk package.
- Tokenization
 - Tokenization is the process of breaking up the given text into units called tokens.
I love Udacity -> [“I”, “love”, “Udacity”]
- Lemmatization
 - In English, words appear in different forms. Look as these two sentence: I have a cat & I have two cats, both sentence talks about cat, but the use in different inflection. Turn these two into the same one by Lemmatization step.
dogs -> dog, cats -> cat

doing -> do, done -> do

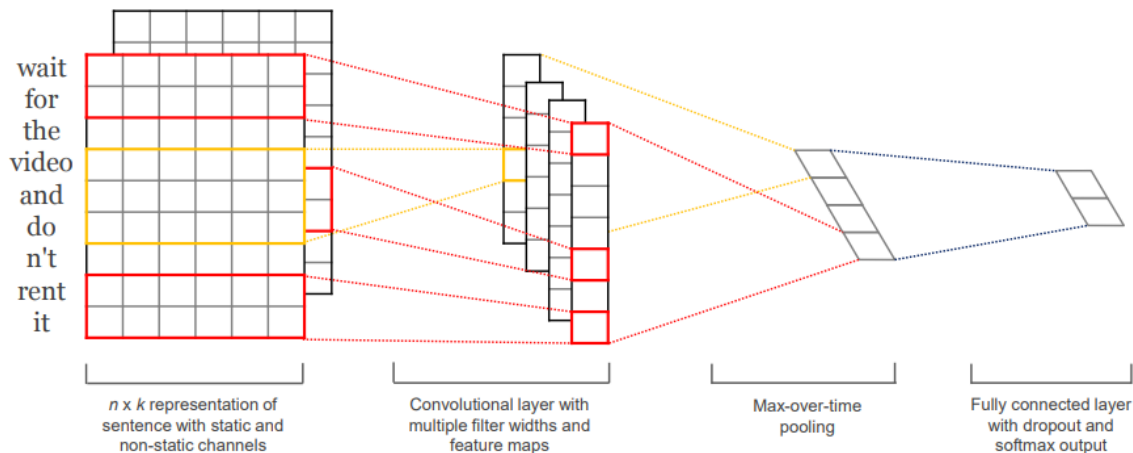
6. Transforming words to Sequences

- Building a dictionary, converting word into vector and every word has same shape of matrix.

```
tokenizer = Tokenizer(num_words=2000) #Create a dictionary of two thousand words and give them index.
tokenizer.fit_on_texts(train_comments)
tr_sequ = tokenizer.texts_to_sequences(train_comments) # convert sentence into matrix
tr_data = sequence.pad_sequences(tr_sequ, maxlen=200)
#We could make the shorter sentences and long sentence have same size matrix
```

7. TextCNN model

- Use Keras framework to build TextCNN model



```
def text_cnn():
    filter_nums = 256
    output_units = 6
    embed_size = 256

    input_layer = Input(shape=(maxlen,), dtype='int32')
    embedding_layer = Embedding(max_features, embed_size, input_length=maxlen,)(input_layer)

    conv_0 = Conv1D(filter_nums, 2, kernel_initializer="normal", padding="same", activation="tanh")(embedding_layer)
    conv_1 = Conv1D(filter_nums, 3, kernel_initializer="normal", padding="same", activation="tanh")(embedding_layer)
    conv_2 = Conv1D(filter_nums, 4, kernel_initializer="normal", padding="same", activation="tanh")(embedding_layer)

    maxpool_0 = GlobalMaxPooling1D()(conv_0)
    maxpool_1 = GlobalMaxPooling1D()(conv_1)
    maxpool_2 = GlobalMaxPooling1D()(conv_2)

    merged_tensor = concatenate([maxpool_0, maxpool_1, maxpool_2])
    output = Dense(units=output_units, activation='sigmoid')(merged_tensor)

    model = Model(inputs=input_layer, outputs=output)
    model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy', auc_roc])
    return model
```

- Embedding_layer: project the words to a defined vector space
- Convolution layer: You are free to increase or decrease the number of layers, and free to adjust the parameters.
- Max pooling layer: receive corresponding the output of convolution layer. Use c
- Concatenate layer: concatenates a list of inputs from Max pooling layer.
- Dense layer: Also called fully connected layer. You are free to increase or decrease the number of layers, and free to adjust the

parameters.

8. Result(accuracy, AUC/ROC score)

Refinement

We constructed three different models architecture below and its performance. Compare the performance of these three models, model 1 and model 2 both have higher performance than model 3. So we will pick one between model 1 and model 2. But model 2 has much more total parameters than model 1, it means it will occupy more resources lead to lower efficiency, when we train model. Incomprehensive survey, I choose model 1 as quasi-final model, and I will make some fine-tuning to this model to make perform better than original one.

Model 1		
Layer	Output Shape	Parameters
InputLayer	(None, 100)	0
Embedding	(None, 100, 256)	512000
Conv1D_1	(None, 100, 128)	65664
Conv1D_2	(None, 100, 128)	98432
Conv1D_3	(None, 100, 128)	131200
GlobalMP1D_1	(None, 128)	0
GlobalMP1D_2	(None, 128)	0
GlobalMP1D_3	(None, 128)	0
Concatenate	(None, 384)	0
Dropout	(None, 384)	0
Dense	(None, 128)	49280
Dense	(None, 32)	4128
Dense	(None, 6)	198
Total Parameters		860,902

Model 1	
Training Accuracy	0.97
Training ROC/AUC Score	0.92
Valid Accuracy	0.93
Valid ROC/AUC Score	0.93

Model 2		
Layer	Output Shape	Parameters
InputLayer	(None, 100)	0
Embedding	(None, 100, 256)	512000
Conv1D_1	(None, 100, 128)	65664
Conv1D_2	(None, 100, 128)	98432
Conv1D_3	(None, 100, 128)	131200
Conv1D_4	(None, 100, 128)	163968
GlobalMP1D_1	(None, 128)	0
GlobalMP1D_2	(None, 128)	0
GlobalMP1D_3	(None, 128)	0
GlobalMP1D_4	(None, 128)	0
Concatenate	(None, 512)	0
Dropout	(None, 384)	0
Dense	(None, 128)	65664
Dense	(None, 32)	4128
Dense	(None, 6)	198
Total Parameters		1,041,254

Model 2	
Training Accuracy	0.97
Training ROC/AUC Score	0.92
Valid Accuracy	0.93
Valid ROC/AUC Score	0.92

Model 3		
Layer	Output Shape	Parameters
InputLayer	(None, 100)	0
Embedding	(None, 100, 256)	512000
Conv1D_1	(None, 100, 128)	65664
Conv1D_2	(None, 100, 128)	98432
GlobalMP1D_1	(None, 128)	0
GlobalMP1D_2	(None, 128)	0
Concatenate	(None, 512)	0
Dropout	(None, 384)	0
Dense	(None, 128)	32896
Dense	(None, 32)	4128
Dense	(None, 6)	198
Total Parameters		713,318

Model 3	
Training Accuracy	0.93
Training ROC/AUC Score	0.89
Valid Accuracy	0.88
Valid ROC/AUC Score	0.88

I found that model performance did not perform well when we constantly increasing convolution-layer, the performance of three convolution-layer structure(model 1) and four convolution-layer structure(model 2) have very close performance. The performance of three convolution-layer structure is better than two convolution-layer structure. But the performance of four convolution-layer structure is almost equal to three convolution-layer structure, so we choose three convolution-layer structure.

Results

Model Evaluation and Validation

I choose model 1 as quasi-final model, and I will make some fine-tuning to this model to make perform better than original one.

Final Model		
Layer	Output Shape	Parameters
InputLayer	(None, 100)	0
Embedding	(None, 100, 256)	512000
Conv1D_1	(None, 100, 128)	131328
Conv1D_2	(None, 100, 128)	196864
Conv1D_3	(None, 100, 128)	262400
GlobalMP1D_1	(None, 128)	0
GlobalMP1D_2	(None, 128)	0
GlobalMP1D_3	(None, 128)	0
Concatenate	(None, 384)	0
Dense	(None, 6)	2310
Total Parameters		809,606

Final Model	
Training Accuracy	0.99
Training ROC/AUC Score	0.99
Valid Accuracy	0.98
Valid ROC/AUC Score	0.98

Justification

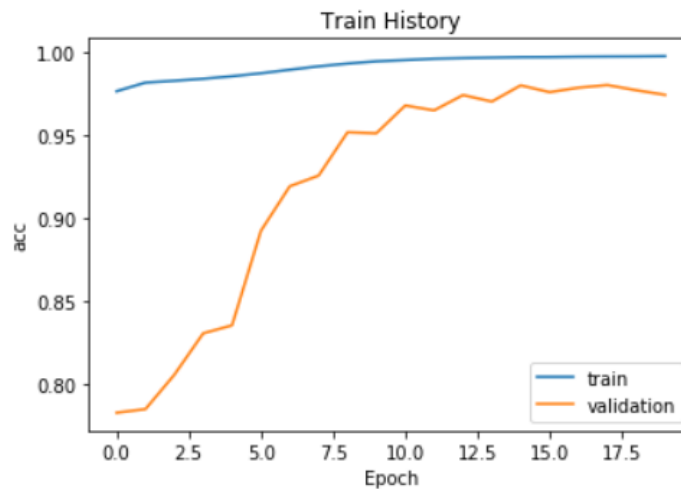
The performance of DecisionTreeClassifier and TextCNN, shown in table below.

DecisionTreeClassifier		TextCNN	
Training Accuracy	0.99	Training Accuracy	0.99
Training ROC/AUC Score	0.85	Training ROC/AUC Score	0.99
Valid Accuracy	0.84	Valid Accuracy	0.98
Valid ROC/AUC Score	0.84	Valid ROC/AUC Score	0.95

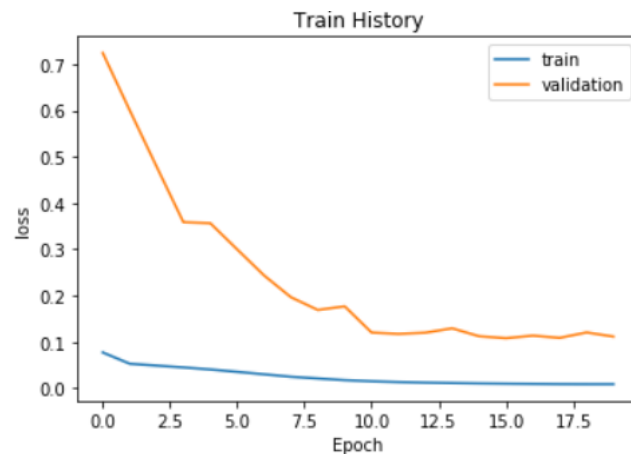
Conclusion

Free-Form Visualization

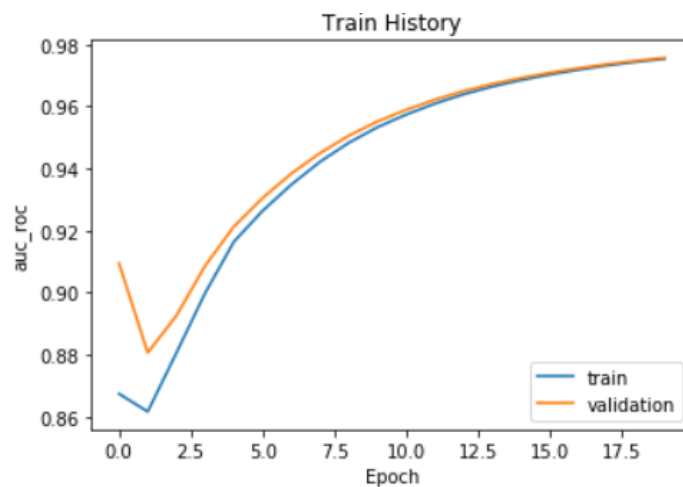
The accuracy of train set and validation set with increasing epoch.



The loss of train set and validation set with increasing epoch.



The AUC/ROC score of train set and validation set with increasing epoch.



Reflection

I spent the most time in cleaning data and I also think that this step is the most important and troublesome in text analysis. Because too much insignificant information can affect the accuracy of the model

Improvement

CNN has good performance in text classification, but there is a deficiency that is not intuitive enough, and the interpretability is not good. The Attention mechanism is a commonly used modeling long-term memory mechanism in the field of natural language processing. It can intuitively give the contribution of each word to the result and does help a lot on catching information in long sequences.

Reference

1. <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge#description>
2. <http://www.aclweb.org/anthology/D14-1181>
3. <http://www.jeyzhang.com/cnn-apply-on-modelling-sentence.html>
4. <https://keras.io/layers/convolutional/>
5. <https://www.kaggle.com/jagangupta/stop-the-s-toxic-comments-eda/notebook>
6. <https://blog.csdn.net/fendouaini/article/details/79919322>