

Toxic Comment Classification

Parth Patel — Vishal Kumar — Amish Ranjan

March 2018

1 Baseline Models

As we have already mentioned in our previous reports that we would be using two models as our baseline model:-

- Logistic Regression
- RNN with LSTM

2 Why we are using these models??

Logistic regression is probably responsible for the majority of industrial classifiers, with the possible exception of naïve Bayes classifiers. Logistic regression is also known as maximum entropy, neural network classification with a **single neuron**, and others. Logistic regression uses unrestricted feature extraction, which allows for arbitrary observations of the situation to be encoded in the classifier. In our situation, we may sometimes have non-standard sentences, which could be better captured using arbitrary observation of Logistic Regression. As aforementioned Logistic Regression can be viewed as single neuron neural network, it is quite obvious that we should try deep neural network with multiple nodes to get better performance. That is why we are trying our second model **RNN with LSTM**.

3 Model Description

- **Logistic Regression:-** Logistic Regression is a classification algorithm. It is used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables. To represent binary / categorical outcome, we use dummy variables. You can also think of logistic regression as a special case of linear regression when the outcome variable is categorical, where we are using log of odds as dependent variable. In simple words, it predicts the probability of occurrence of an event by fitting data to a logit function. It is easy to minimize loss function:-

$$l(\theta) = \sum_{i=1}^n y^i \log(h(x^i)) - (1-y^i)(1-h(x^i))$$

If there is a little deviation from actual function, there is exponential higher loss for that.

- **RNN with LSTM:-**Unlike the conventional translation models, where only a finite window of previous words would be considered for conditioning the language model, Recurrent Neural Networks (RNN) are capable of conditioning the model on all previous words in the corpus. But RNN has Vanishing Gradient Gradient Explosion Problems, which are solved by LSTM. That is why we are using RNN with LSTM.

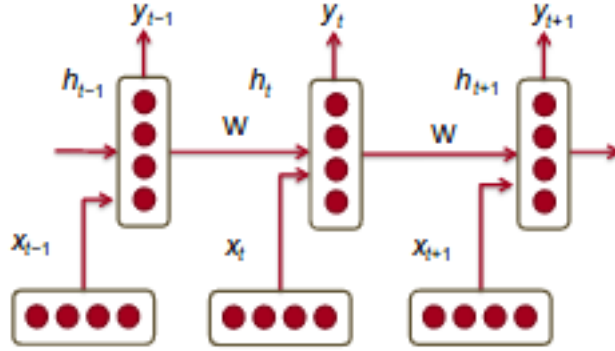


Figure 1: A recurrent neurall network with three time steps shown

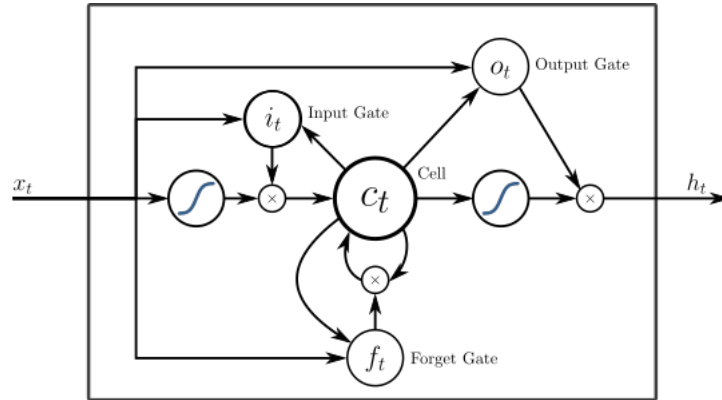


Figure 2: An LSTM with all gates

LSTM model which introduces a new structure called a memory cell (see Figure 1 below). A memory cell is composed of four main elements: an input gate, a

neuron with a self-recurrent connection (a connection to itself), a forget gate and an output gate. The self-recurrent connection has a weight of 1.0 and ensures that, barring any outside interference, the state of a memory cell can remain constant from one timestep to another. The gates serve to modulate the interactions between the memory cell itself and its environment. The input gate can allow incoming signal to alter the state of the memory cell or block it. On the other hand, the output gate can allow the state of the memory cell to have an effect on other neurons or prevent it. Finally, the forget gate can modulate the memory cell's self-recurrent connection, allowing the cell to remember or forget its previous state, as needed.

4 Data Collection

We have referred to Kaggle Competition, for this project. There is a well defined and structured dataset available. The dataset consists of Wikipedia comment made available by Google-Jigsaw in public domain. The training set have:-

id	comment_text	toxic	\
0000997932d777bf	Explanation\nWhy the edits made under my usern...	0	
000103f0d9cfb60f	D'aww! He matches this background colour I'm s...	0	
000113f07ec002fd	Hey man, I'm really not trying to edit war. It...	0	
0001b41b1c6bb37e	"\nMore\nI can't make any real suggestions on ...	0	
0001d958c54c6e35	You, sir, are my hero. Any chance you remember...	0	

severe_toxic	obscene	threat	insult	identity_hate
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

Figure 3: An example of test set

- id
- comment line
- toxic(0,1)
- severe_toxic(0,1)
- obscene(0,1)
- threat(0,1)
- insult(0,1)
- identity_hate(0,1)

5 Data Pre-processing

The data is given in English language, we have parsed it word by word to calculate tf-idf(with n-gram). We have also pre-processed the non-standard english like (a\$\$hole, l o s e r) to make it useful for our training purpose, otherwise it was being discarded.

6 Result

We are using 2-3rd of the data as training and the rest for test.

- Logistic Regression :-

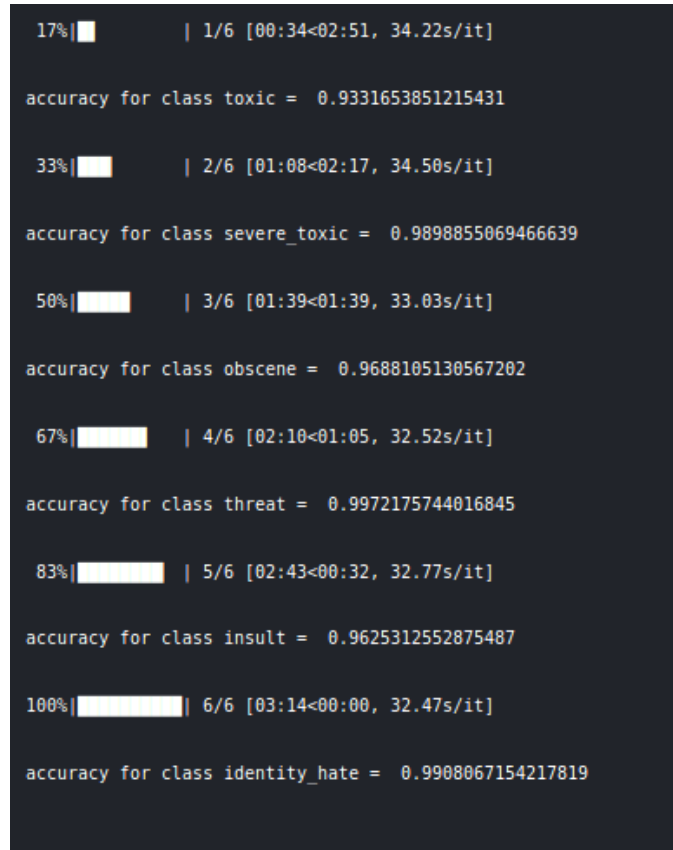


Figure 4: Results for logistic regression

- Multi Layer Perceptron :- We have used five hidden layers for the multi-layer perceptron.

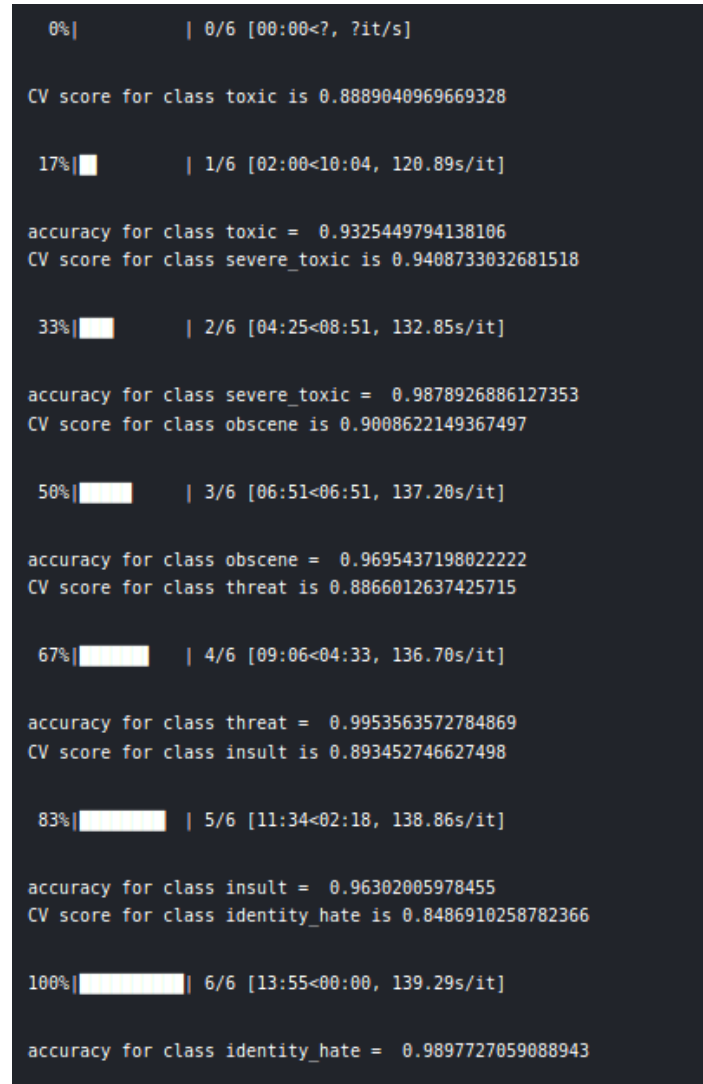


Figure 5: Results for Multi Layer Perceptron

- LSTM In Deep Neural Network :- We have used one layer of LSTM with 300 LSTM cells and one dense layer with one node for output.

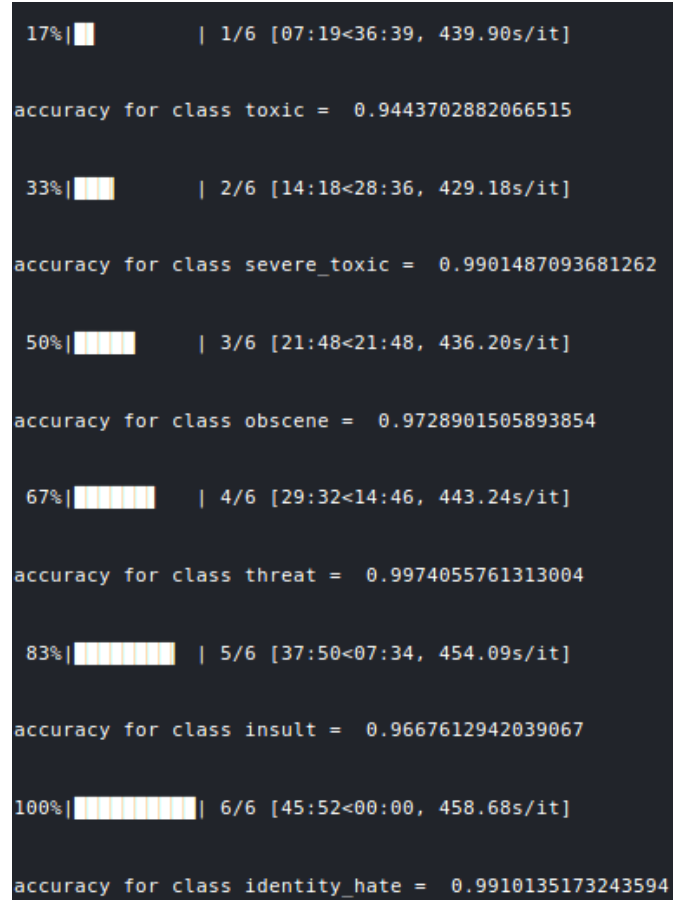


Figure 6: Results for LSTM

- Deep Neural Network :- We have used one Dense layer of 100 nodes with 'relu' activation, next layer is again dense layer with 10 nodes with same 'relu' activation and last layer is again dense layer with one node for output with sigmoid activation.

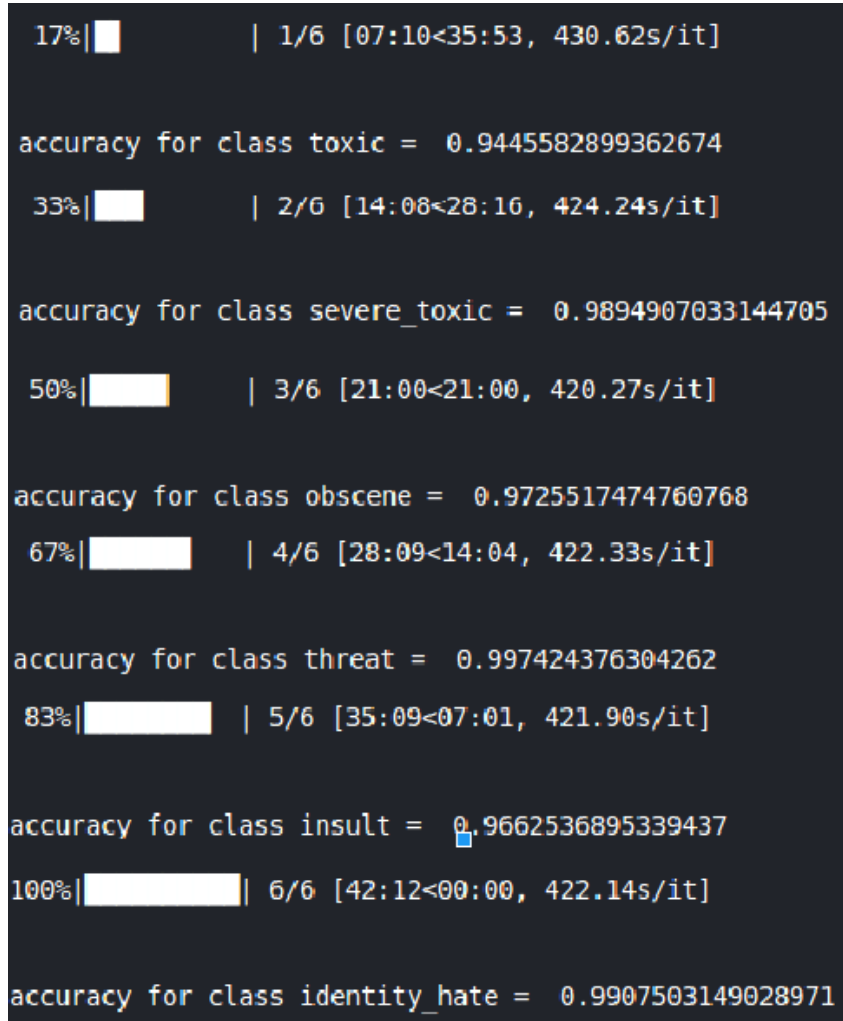


Figure 7: Results for Deep Neural Network