



# Toxic Comment Classification

CMPT 413 Final Project

Group : outputerror

Group Member: Ao Tang, Shaoqiang Zou



# Motivation

- The internet has enabled people to communicate and learn from each other
- Online hate, abuse, toxicity shuns people from opening up and taking full advantage of the opportunities that online communication provides
- Our goal is to use machine learning techniques to classify toxicity from the six types of toxicities



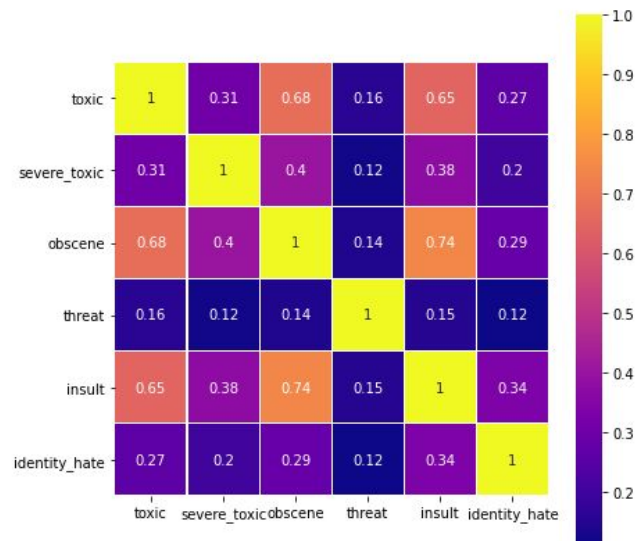
# Dataset

- Comments from Wikipedia's talk page
  - Train: 150k comments; Test: 60k comments
- Six types of toxicities are all boolean labels (0 or 1)

|          |                |                 |
|----------|----------------|-----------------|
| ➤ toxic  | ➤ severe_toxic | ➤ obscene       |
| ➤ threat | ➤ insult       | ➤ identity_hate |

# Dateset

|   | 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             |

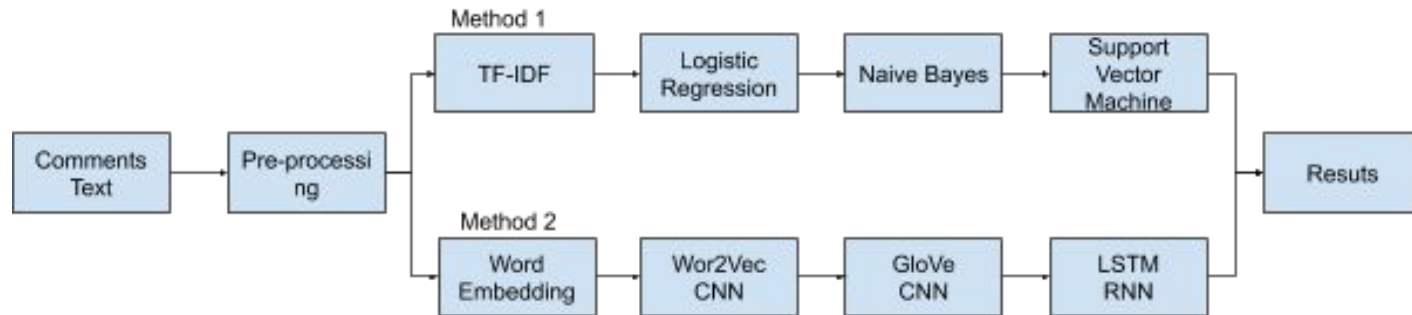




# Data Pre-processing

1. Translate Non-English word to English word use Google Translate API
2. Convert all character to lowercase
3. Replace emoticons with word
4. Replace Date, Phone Number, and Website Links
5. Remove all numbers and punctuations
6. Normalize repeating characters
7. Remove stopwords

# Approach



## Method 1: TF-IDF

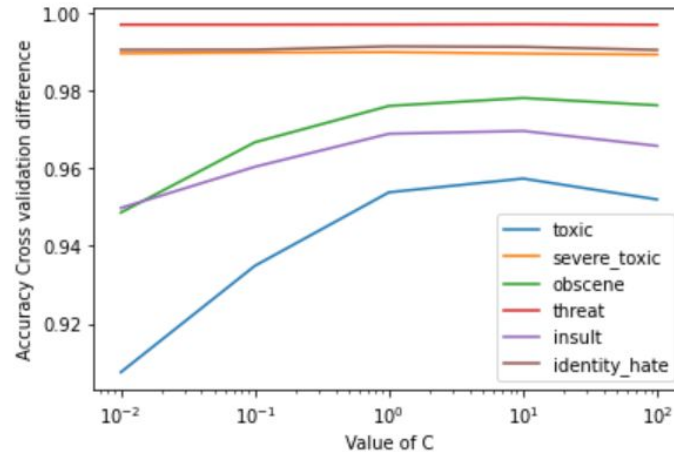
- Comments in document is transform into Term Frequency/ Inverse document matrix
- Every unique word is considered as feature
- TF-IDF score can be fed to logistic regression, Naive Bayes, Support Vector Machine

```
print(tf_vec.get_feature_names())
```

['a', 'aaba', 'aabove', 'aacd', 'aachen', 'aachi', 'aadd', 'aadm', 'affect', 'afia', 'aflight', 'aafs',  
 'aagadu', 'aages', 'aagr', 'aagh', 'aagin', 'aah', 'aahank', 'aahh', 'aahs', 'aai', 'aaiha', 'aajacksoniv',  
 'aajonus', 'aakash', 'aake', 'aalborg', 'aalertbot', 'aalexa', 'aaliya', 'aaliyah', 'aaliyahremembered', 'aal',  
 'aalst', 'aamir', 'aand', 'aanda', 'aandah', 'aandc', 'aande', 'aandm', 'aandr', 'aandw', 'aang', 'aanny',  
 'wwheerree', 'aany', 'aao', 'aapropritate', 'aarabs', 'aaru', 'aardman', 'aardsma', 'aardsman', 'aardvark',  
 'aardvarks', 'aare', 'aarem', 'aargh', 'aaronhorn', 'aare', 'aaroamal', 'aaroohi', 'aaron', 'aaroncric', 'aa',  
 'aronic', 'aaronsw', 'aaround', 'aarp', 'aarrgh', 'aarrow', 'aasc', 'aave', 'aviskoo', 'avishkarh', 'aaw', 'a',  
 'aaw', 'aaye', 'aayege', 'abacha', 'abacination', 'aback', 'aback', 'abacus', 'abad', 'abadan', 'abaddon', 'aba',  
 'de', 'abagnale', 'abalessa', 'abali', 'abangani', 'abandon', 'abandoned', 'abandoned', 'abandoning', 'abandonm',  
 'ent', 'abandonou', 'abandons', 'abanes', 'abang', 'abanteart', 'abaranger', 'abarenoh', 'abase', 'abased', 'a',  
 'bassids', 'abatayo', 'abate', 'abaya', 'abba', 'abbas', 'abbasgulu', 'abbasid', 'abbasids', 'abbass', 'abbasse',  
 'd', 'abbassi', 'abbassid', 'abbastanza', 'abbau', 'abbe', 'abberations', 'abberline', 'abberant', 'abbes', 'a',  
 'bbee', 'abbeys', 'abbi', 'abbie', 'abhiit', 'abbot', 'abbot', 'abbotabad', 'abbotandcostello', 'abbotstfor',  
 'd', 'abbr', 'abbrasive', 'abbreve', 'abbreviated', 'abbreviate', 'abbreviated', 'abbreviateds', 'abbreviates',  
 'abbreviating', 'abbreviation', 'abbreviations', 'abbrevs', 'abbreviations', 'abbron', 'abby', 'abbytes', 'abb',  
 'ythead', 'abbywinters', 'abce', 'abce', 'abcedere', 'abcedere', 'abcmmonster', 'abcnws', 'abdali', 'abdalla',  
 'h', 'abdallar', 'abdaly', 'abdalyar', 'abdel', 'abdellaziz', 'abdellhasen', 'abdellkarim', 'abdi', 'abdicate', 'a',  
 'bication', 'abdielcolberg', 'abdillah', 'abdin', 'abdnor', 'abdolhassen', 'abdolmalek', 'abdomen', 'abdomina',  
 'l', 'abdozy', 'abdozora', 'abdou', 'abdou', 'abducted', 'abducting', 'abduction', 'abductions', 'abductive',  
 'abdu', 'abdulaziz', 'abdulkhaki', 'abdulhamid', 'abdulkadir', 'abdulla', 'abdullah', 'abdullahi', 'abdulrah',

# Method 1: Logistic Regression

- Six models for each type of toxicity
- Use cross validation to tuning parameter







# Method 1: Logistic Regression

- Used Accuracy for measuring model performance on test data set.
- Our results from logistic are following:

| Type of Toxicity | Performance(Accuracy) |
|------------------|-----------------------|
| Toxic            | 0.9574                |
| Sever_Toxic      | 0.9900                |
| Obscene          | 0.9782                |
| Threat           | 0.9972                |
| Insult           | 0.9697                |
| Identity_Hate    | 0.9915                |

- 
- We also worked on other models, the average accuracy score are following:

| Models                       | Performance(Accuracy) |
|------------------------------|-----------------------|
| Logistic Regression          | 0.9722                |
| Naive Bayes                  | 0.9708                |
| Support Vector Machine (SVM) | 0.9721                |

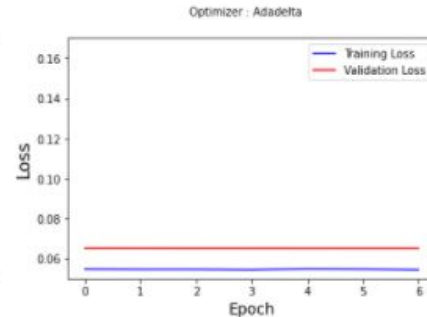
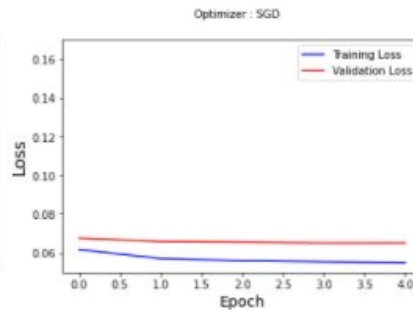
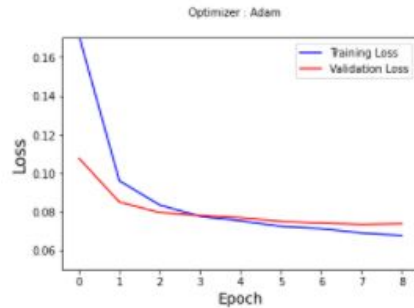


## Method 2: Word Embedding

- **Word2vec - Convolution Neural Nets**
  - Pre-trained model from Google
  - Contains 300 - dimensional vectors for 3 million words
- **GloVe - Convolution Neural Nets**
  - Pre-trained model from Stanford
  - Contains 300 - dimensional vectors for 400,000 words
- **Keras word embedding - Recurrent Neural Nets**

# Experiments Performed

- Early stopping using validation set
- Tried different optimisers (Adam, SGD, and AdaDelta)





# Model Performance

Average accuracy score:

| Model | Embedding               | Validation | Test   |
|-------|-------------------------|------------|--------|
| CNN   | word2vec                | 0.9789     | 0.9732 |
| CNN   | GloVe                   | 0.9789     | 0.9731 |
| RNN   | keras<br>word-embedding | 0.9799     | 0.9722 |



# Conclusion

- TF-IDF: Logistics Regression gives best performance.
- Word Embedding: Word2vec gives best performance.
- Overall, our word embedding method outperformed TF-IDF method.



# Dataset Credits

Jigsaw/Conversation AI. 2017.

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