

Toxic Message Detection

Machine learning

HSE University

2023-06-13

Primary preprocessing

- Tokenize
- Drop stop words

(file: preprocess.py)

Frequency preprocessing: Bayesian analysis

How many times is this word encountered in toxic and regular messages?

⇒ how much does the given word affect the probability, that the message is toxic?

Frequency dictionary contains 150 words with the highest effect

Frequency preprocessing: formula

- T - total number of toxic messages
- t - number of times this word appears in toxic messages
- T^c and t^c are their respective complements

$$effect = \left(\frac{t^c}{T^c} - \frac{t}{T} \right) \frac{T + T^c}{t + t^c}$$

(file: preprocess_frequency.py)

RNN Preprocessing: Words

- glove.840B.300d words semantic database
- Turn each tweet into a list of vectors of length 300

(file: preprocess_rnn.py)

RNN Preprocessing: Training

- Supply encoder initial state and, sequentially, input tokens
- Obtain the final state
- Check every decoder's output against the expected input
- Update weights

The output is a 100d vector. (file: preprocess_rnn.py, rnn.py)

Fit models: Frequency Dictionary

- Input: vector of 150 floats. Each float represents how many the given message contains words from the 150 most influential words.
- Output: class 0 or 1

Fit models: RNN

- Input: vector of 100 floats, which correspond to the semantics of the message
- Output: class 0 or 1

Models fitting and comparison

Model	Train error	Test error
Freq SVM	0.9766	0.9568
RNN AdaBoost	0.9694	0.9614