# Predicting Closed Stack Overflow Questions

**Levi Franklin**
Department of Computer Science
Stanford University
Stanford, CA
*lefrankl@stanford.edu*

## Abstract

Stack Overflow is a website that is used by software developers extensively to ask and answer questions about software development. Each day they reportedly receive over 6000 questions. Due to this high volume of questions it is difficult and time consuming to analyze them all and determine which should be left open and which should be closed for a variety of reasons.

This project aims to provide a method to assist in detecting what questions should be closed and which should remain open. Looking at this problem as a sentence classification problem one can use deep learning and neural network models to train a system that will identify good candidates for closing. With the help of a large dataset provided by Kaggle and Stack Overflow, one can use algorithms that are effective in tasks such as sentiment classification and see how they apply to this task.

To evaluate the success of these models, they are compared to several benchmarks as well as to other competitors in the competition. The multiclass log loss is used as it was in the competition to facilitate this comparison. Different models and sets of word vectors are used and their log losses are compared to determine which system works best.

In addition to the natural language processing that is done on the text of the question, there are some numeric features which are utilized such as the question asker's reputation. Similar tests will be done to see if incorporating these features in the model provides better results. Overall this project aims to show that a good estimate of what questions should be closed can be given by a neural network based prediction model.

## 1    Introduction

Stack Overflow is a website that is used by developers throughout the world to ask and answer questions that they may have related to software development. Among the reportedly 6000 questions Stack Overflow receives every weekday, many are inadequate to make it onto the site. These questions can be closed for being off topic, not constructive, not a real question, or too localized. With this many questions it is difficult to inspect each one so an automated process to assist in determining what should be closed is very helpful.

The was recently a Kaggle competition to solve this very issue which includes a dataset of questions, metadata about each question, and whether it was left open, or the reason it was closed. This dataset provides a good set of training and testing data to build a neural network

43  to try to create an automated system for making the determination on if a question should be
44  closed.

45  The first task to consider is how to model the actual words that make up each of the questions
46  in our neural networks. Word vectors are used to represent the words within each question and
47  consist of numeric vectors of some dimension (often 50-100). These word vectors are a
48  representation of how words appear in relation to each other in a large corpus of documents.
49  This project attempts using a pretrained set of word vectors that is well regarded in the field,
50  the GloVe vectors. Additionally, a new set of vectors is trained using the Stack Overflow
51  dataset which incorporates the technical stack overflow terms that may not be in the GloVe
52  set. These two sets of word vectors are compared to see which has better performance.

53  There are many neural network models that could be applied to this problem. Primarily for
54  this project there are two main neural network models that are considered. The first is a system
55  that averages these word vectors and uses those as the input to a traditional multi-level neural
56  network is used. This has some shortcomings such as not taking into account sentence
57  structure. Secondly a convolutional neural network model is used that passes over the words
58  in each question to consider them in groups. Each of these filters, or convolutions, outputs
59  some estimate that is then pooled together. The result of this is used to make a prediction on
60  what class that input belongs to. This does a better job of including sentence structure but is
61  more complicated and time consuming to run.

62  The other portion to consider is how to incorporate numeric features in addition to these word
63  vectors. Each question includes some numeric features such as the user's reputation at the time
64  of asking or how many open questions they have at the time of asking. The main method of
65  including these features is taking them on to the end of each word vector so that they can be
66  included in the training. Training with these features is compared to training without to see
67  which is more helpful.

68  In order to gauge the performance of these models the log loss metric is used. This is the
69  metric that was used during the Kaggle competition which will allow comparison with other
70  competitors. Additionally, the loss of these models will be compared to some baselines. This
71  will include making a uniform prediction of each class as the most naïve baseline, as well as
72  one based purely on the class distribution in our training data. Lastly Stack Overflow has
73  provided a machine learning model they created to compare against other submissions.
74  Overall this project aims to beat both more rudimentary baselines and be near many of the
75  other Kaggle submissions in performance.

76

## 2    Background/Related Work

78  This problem can be treated primarily as a classification problem where the main goal is to
79  take in one stack overflow question and the metadata associated with it, and classify it as one
80  of 5 categories: open, off topic, not constructive, not a real question, or too localized. This
81  problem is a familiar topic in natural language processing called sentence classification. It has
82  been studied before and many people have proposed solutions. A common task that is used as
83  an example in these studies is doing sentiment analysis of things like movie reviews [1]-[4].

84  The first aspect in these sort of problems is representing the words in these sentences as word
85  vectors. These word vectors are generally generated by creating a matrix of what words are
86  near each other in a large corpus. Usually this corpus consists of things like Wikipedia. One
87  popular word vector representation is the GloVe system which creates a coooccurence matrix
88  and then reduces it using techniques like SVD [5]-[6]. Another model is called word2vec or
89  the skipgram model and learns word vectors on the fly using a predictive model [7]. One can
90  use vectors trained in either way for a neural network model.

91  Once a method or representing the questions has been determined, it is prudent to investigate
92  what models will best make the associated classification. One easier yet surprisingly powerful
93  solution is to just average the word vectors that are used in the sentences. This does not hold
94  the structure of the sentence but can capture some meaning. Another technique makes use of
95  recursive neural networks which feeds the sentences into a set of similar network levels that
96  share a single weight matrix and feed into the next one, capturing the result of all the words
97  before it [8]. This model suffers from a vanishing gradient issue where it forgets about older

98     words as time goes on and the gradient becomes weaker.

99    Two other techniques help lessen this gradient problem by using as building blocks of their
100   network a layer that has some "memory". Gated Feedback Recurrent (GRUs [9]-[10]) and
101   Long Short-Term Memory (LSTMs [11]) neural networks keep track of more information from
102   previous parts of the sentence.

103   Another system that takes into account structure is convolutional neural networks [3]. These
104   networks take a filter and do a convolution over the sentence. It looks at some number of
105   words and passes along the sentence making predictions based on those words. This can
106   happen multiple times and in the end the results are pooled to make a final prediction. Overall
107   these technique provide several different ways that sentences can be processed and classified
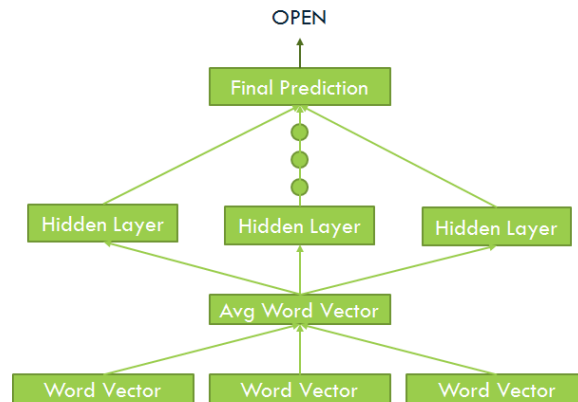108   based on a set of training data.

109

## 110   3     Approach

111   The approach that is taken in this paper explores two different models and compares them with
112   various input types. Additionally, it makes use of two different sets of word vectors and
113   compares them to gauge which is more applicable. These different techniques are trained
114   against a large dataset of questions provided by Stack Overflow and are then tested on a
115   separate subset of those questions. The baselines are also computed on that training set and
116   are used to compare these values to the other Kaggle competitors.

117   The two kinds of word vectors that are considered are the GloVe vectors that are mentioned
118   above and ones that are trained specifically on the corpus of training data that we are using.
119   The GloVe vectors are pre trained and are known to do a good job at expressing the ways that
120   word correlate with one another. However, they are trained on a corpus of Wikipedia and
121   Gigaword [12] so they may be missing some technical terminology that Stack Overflow uses.
122   In an effort to overcome this a large set of Stack Overflow questions is also used to generate
123   another set of word vectors that may be better suited to this task.

124   Next the first model that is tested should be discussed. This model is an averaging of all the
125   word vectors in each sentence. This resulting word vector is then fed into a neural network
126   and is trained against the dataset of questions. This neural network takes in the average word
127   vector, runs it through some number of hidden layers which multiply it by various weights
128   which are learned during training. The final layer combines all of these into a final prediction
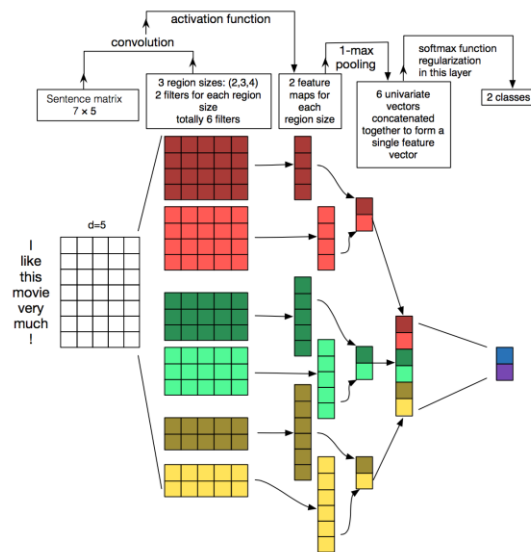129   of open or closed.

130



131                          Figure 1. Averaging Word Vectors

132   The next model that is used is a convolution neural network model. This type of model takes
133   filters of various sizes and applies them over the word vectors. These filters are passed over
134   the words in the sentences in sets that are the size of each filter. Then these filters outputs are
135   pooled together and used to make a prediction much like in the model before. In this model,
136   the actual structure of the sentences matters and influences the final results.

137

Figure 2. Convolutional Neural Network [13]

The other main consideration is how to incorporate the numeric features in each of these models. As was mentioned before, each question also includes the reputation of the user and the number of other open questions they have. This information can be very helpful in determining if a question should be open as someone with high reputation is much more likely to submit a question correctly. The technique that is used here is to include those two values in the end of the word vectors. For averaging word vectors it is done after they are averaged and for the convolutional network they are added to every vector. This allows them to be included in the neural computations just like any other feature.

Lastly some discussion of the metric that is used to evaluate these models is important. The Kaggle competition used the mutli class log loss. As will be seen later on the dataset that is used has a highly biased distribution towards the open category. This means most all models will achieve high accuracy as well as very similar precision and recall so these metrics are not overly helpful. In order to compare these models to the baselines and the other competitors the log loss is used here as well. The equation for the log loss multiplies the log of the probability for the correct class of each input and averages them all.

$$\text{Log Loss:} -\frac{1}{N}\sum_i y_i \log(\overline{y_i})$$

## 4    Experiment

### 4.1    Dataset

The dataset that was used is a very large collection of Stack Overflow questions coming from Kaggle and Stack Overflow. In total it has around 3.4 million questions. The dataset distribution is heavily weighted towards the open class. Stack Overflow has stated that around 6% of their questions are close with even less being marked as closed in the training data. This data was split up as 60% training, 20% dev, and 20% final testing.
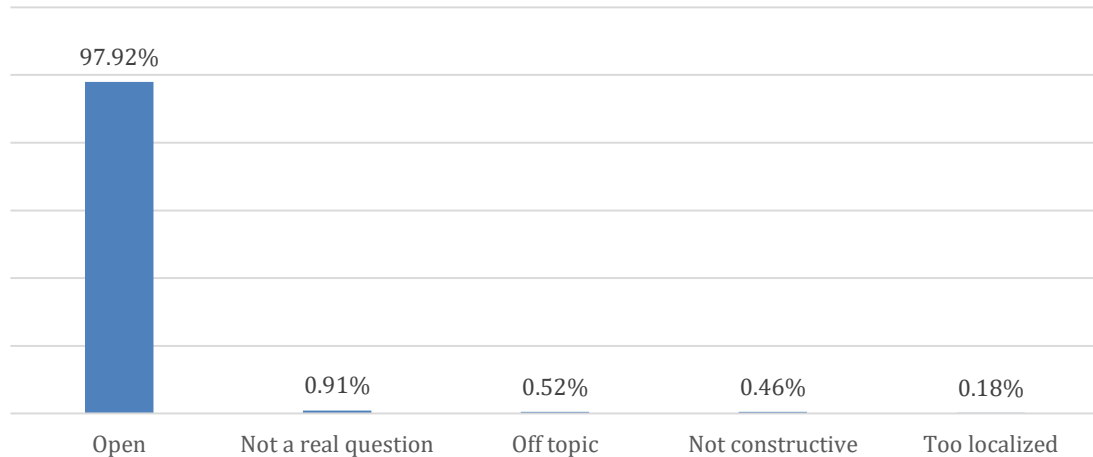
| | | |
|---|---|---|
| **Training** | 2,022,318 | 60% |
| **Development** | 674,106 | 20% |
| **Testing** | 674,106 | 20% |

Figure 3. Distribution of 3,370,530 questions in dataset

## Distribution of Class Labels

97.92%

0.91%    0.52%    0.46%    0.18%

Open    Not a real question    Off topic    Not constructive    Too localized

Figure 4. Distribution of class labels

The dataset includes the following features that are useful in training a neural network:

- Question Title
- Question Text
- Submitter's stack overflow reputation
- Number of questions submitter has open at time of submission

Here is a sample entry for a question that was marked as closed:

**Users reputation:** 32

**Users other open question count:** 0

**Title:** List of all .txt file

**Text:** I want to write a program that give a path in my system and goes to that path and search in that path and sub-directory of path and list all of .txt file. please help me thanks .

**Label**: not a real question

### 4.2    Results

To begin the values of log loss on the test set for each of the baselines was collected. As was mentioned before the uniform baselines predicts evenly across all classes and provides a very bottom baseline. Next is a baseline that predicts the same distribution as the dataset distribution that was mentioned above. This is a better baseline in that it predicts open heavily and is almost always right. The goal is to improve upon this baseline. The last baseline is the one based on the model that was given by the Stack Overflow team as their baseline for predicting. This improves upon the frequency distribution baseline substantially and is a good goal for these models to achieve.

The first experiment was to train the averaging model with both GloVe word vectors and the Stack Overflow trained vectors, then compare the final log loss with the baselines and the other Kaggle competitors. Unfortunately, the correct labels for the final submission in the Kaggle competition were not made available so this comparison can only be made in comparison to these three baselines. Secondly, the same experiment is done using the convolutional model. On this experiment however, the GloVe vectors were excluded as it was seen that the Stack Overflow trained word vectors performed better. There was also some work done to try resampling the dataset to compensate for the uneven class distribution, and then readjusting the predictions based on their correct distribution in Stack

200    Overflow questions. However, this did not prove to help and was abandoned.

201

| Method | Log Loss |
|---|---|
| Uniform Prediction Baseline | 1.609 |
| **Averaging Stack Overflow Trained Word Vectors** | **0.205** |
| **Averaging GloVE Word Vectors + Numeric Features** | **0.200** |
| **Convolution Stack Overflow Trained Word Vectors + Numeric Features** | **0.212** |
| **Convolution Stack Overflow Trained Word Vectors** | **0.190** |
| Frequency Based Prediction Baseline | 0.173 |
| **Averaging Stack Overflow Trained Word Vectors + Numeric Features** | **0.152** |
| Stack Overflow Model Prediction Baseline | 0.094 |

202    Additionally, here is the confusion matrix for the Averaging Stack Overflow Trained Word
203    Vectors + Numeric Features results:

| | Not a real question | Not constructive | Off topic | Open | Too Localized |
|---|---|---|---|---|---|
| Not a real question | 4 | 8 | 20 | 9856 | 0 |
| Not constructive | 3 | 63 | 21 | 3438 | 0 |
| Off topic | 2 | 10 | 97 | 4498 | 0 |
| Open | 6 | 54 | 172 | 653735 | 1 |
| Too Localized | 0 | 1 | 0 | 2114 | 1 |

204    And the precision and recall for each entry:

| Label | Precision | Recall |
|---|---|---|
| not a real question | 0.2667 | 0.0004 |
| not constructive | 0.4632 | 0.0179 |
| off topic | 0.3129 | 0.0211 |
| open | 0.9705 | 0.9996 |
| too localized | 0.5000 | 0.0005 |

205    Here is an example of correctly classified document as "not a question":

206        **Title:** MySQL Table of UK Area Codes and Names
207        **Text:** I need database of Uk cities with their post codes. Any help is appreciated.
208        Thanks for your time..

209    # 4    Conclusions

210    Overall this has been a fairly successful experiment. Though many of the models did not achieve
211    outstanding results, the best performing one did manage to beat both of the starting baselines. This
212    shows that a neural network can indeed be used to help facilitate predicting what questions on

213 Stack Overflow should be closed. As can be seen by the confusion matrix, it was able to correctly
214 predict questions in every category.
215
216 Using this information this models performance can also be compared to those of the other Kaggle
217 competitors. Though it cannot be compared directly due to not having the labels on the test set
218 they used, they can be roughly compared due to the benchmarks. Based on those, the best model
219 here would come around 110th out of 160 submitters. Assuming Kaggle competitors are usually
220 well versed in machine learning, this seems to be a pretty good result.
221
222 Interesting questions are what caused the other models to perform not as well. This could be due
223 to bugs in the models or not achieving the best hyperparameters. Each of them have various things
224 that can be tweaked such as number of layers, number of filters, and many others. Though lots of
225 various parameters were used to find the best one, it is possible the best parameters were not
226 achieved. Other future work could entail using other features from stack overflow such as the
227 number of questions the user has had cancelled before. Additionally, the other models that were
228 mentioned during the related work could provide better results.

229 **References**

230 [1] Andrew L. Maas and Andrew Y. Ng (2010) **A Probabilistic Model for Semantic Word Vectors**

231 [2] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher
232 Potts (2011) **Learning Word Vectors for Sentiment Analysis**.

233 [3] Yoon Kim (2014) **Convolutional Neural Networks for Sentence Classification**

234 [4] Quoc Le, Tomas Mikolov (2014) **Distributed Representations of Sentences and Documents**

235 [5] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. (2014). **GloVe: Global Vectors**
236 **for Word Representation.**

237 [6] Eric H. Huang, Richard Socher∗ , Christopher D. Manning, Andrew Y. Ng (2015) **Improving**
238 **Word Representations via Global Context and Multiple Word Prototypes**

239 [7] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, Jeffrey Dean (2013) **Distributed**
240 **Representations of Words and Phrases and their Compositionality**

241 [8] Toma´s Mikolov,  Martin Karafiat, Luka´s Burget, Jan "Honza" Cernock, Sanjeev Khudanpur
242 (2010) **Recurrent neural network based language model**

243 [9] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, Yoshua Bengio (2014) **Empirical Evaluation**
244 **of Gated Recurrent Neural Networks on Sequence Modeling**

245 [10] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, Yoshua Bengio (2015) **Gated Feedback**
246 **Recurrent Neural Networks**

247 [11] Sepp Hochreiter, Jurgen Schmidhuber (1999) **Long Short-Term Memory**

248 [12]   Parker, Robert, et al. English Gigaword Fifth Edition LDC2011T07. DVD.
249 **Philadelphia: Linguistic Data Consortium**, 2011.

250 [13]  http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/
251 (2015) **UNDERSTANDING CONVOLUTIONAL NEURAL NETWORKS FOR NLP**