

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Classifying responses on online discussion forums

Aaron Abajian

aabajian@stanford.edu

Abstract

There are online discussions forums covering an endless variety of topics. Each forum has many topic-specific questions with expert answers. However, the best answer is often buried somewhere in the discussion thread. Given a question on an online forum, we focus on the problem of classifying each response to the question as an *answer* or a *non-answer*. Where the *answer* label indicates if the response made a reasonable attempt at answering the question (even if the answer is incorrect).

We train a recursive neural network (RNN) using 59,205 question-response pairs from ten online forums with true labels obtained through Amazon Mechanical Turk. Our model reads each question-response pair from left-to-right similar to a next-word prediction RNN. The model minimizes the number of incorrectly-classified question-response pairs. After hyperparameter tuning, the model achieved a 76.3% classification accuracy. We qualitatively validate the model by showing common *non-answer* responses such as “thanks”, “bump”, and “nevermind” result in strong neuronal activation.

Our work builds off of previous work done by Aaron Abajian as part of CS 229 and CS 221. In those classes, a logistic regression (LR) classifier and a support-vector machine (SVM) classifier were trained using hand-selected features. The SVM and LR classifiers outperformed the RNN due to hand-crafted features. We suggest ways for improving the RNN to match their performance.

1 Introduction

Online discussion forums have existed since the early days of the Internet. A forum consists of threads with one or more posts. An original poster (OP) creates a thread by posing a question pertinent to the forum’s topic. The OP can expect accurate answers because niche forums attract experts in their field.

One problem with forums is that good answers are often buried somewhere in the discussion thread. Traditional forums do not provide functionality to re-order responses based on *correctness*, they merely list responses chronologically. Moreover, users are free to respond with *non-answers* such as: “I have this same problem!”, “Can you provide more information?”, “No idea, good luck.”, and “Thanks, that solved it.” None of these responses constitute *answers* to the OP’s question. New visitors with the same question must invariably sift through non-answers to find answers. This often leads to the creation of a new thread with the same question and the ubiquitous non-answer response, “This has been asked many times, why not use the search function next time?”

Question and Answer (Q&A) sites are a recent trend that attempt to resolve these issues. Sites such as StackOverflow, Answers.com, Yahoo! Answers, and Quora emphasize *answer*-only responses to questions that are up-voteable by users. *Non-answers*, such as the examples above, are delegated to comment sections or down-voted from view.

054 It would be beneficial if responses on traditional forums could be up-voted automatically. In this
055 project, we address a relaxed version of this problem. We seek to distinguish *answers* from *non-*
056 *answers*. We define an *answer* to be any response that makes a reasonable attempt at answering the
057 question.

058 We train a recurrent neural network (RNN) that assigns a probability to each question-response
059 pair that indicates how likely the response is to answer the question. Responses with less than 0.5
060 probability are classified as *non-answers*. Responses may be re-ordered based on their probability
061 assignments so that non-answers are moved to the bottom of the thread. A forum search engine that
062 excludes non-answers increases the likelihood that new visitors will find the answer to their question
063 in a reasonable amount of time.

064 Note that this work builds off of previous work performed by Aaron Abajian in CS 229 and CS 221.
065 In those projects, a support-vector classifier and a logistic regression classifier were trained on the
066 same classification task. We compare performance of the RNN to those models in the discussion
067 section.
068
069
070
071

072 **2 Background / Related Work**

073
074

075 Question-answering (Q&A) systems have received considerable attention in the past two decades. A
076 variety of models have been trained using decision trees [1], support vectors machines (SVMs) [6]
077 [7], and recurrent neural networks [4]. The majority of these models are trained using large databases
078 of facts and factoids (small snippets of knowledge). Given a question, these models compute a
079 probability distribution over all known facts and return those facts with the highest probabilities.

080 The decision tree and SVM models typically relay upon n -gram feature vectors and cosine simi-
081 larity. They may also include special features such as date indicators, parts-of-speech tagging, and
082 named entity tagging. Their performance is limited because they cannot take into consideration
083 long-distance interactions such as questions with multiple sentences and deep relationships between
084 words, phrases and their context.

085 The neural network models, in particular [4], have been more successful than traditional n -gram
086 based word-matching schemes because they are able to take into consideration the “syntactic glue”
087 that gives meaning to a question and its context. Higher-dimensional learned word and phrase
088 vectors have been shown to capture deep syntactical relationships [5]. For example, Pennington, et.
089 al. show the close proximity to the word vectors for “queen - woman” and “king - man”. Perhaps
090 more poignant for question-answering is the ability of learned word vectors to identify new words
091 that mean approximately the same thing. For example, “litoria” and “leptodactylidae” are both
092 appropriate responses to the question, “What are some frog genera?” Their learned word vectors are
093 both similar to “frog”, a fact that is quite well hidden in n -gram only features.

094 Fact-based Q&A systems are inherently limited by their corpus of training data. This is a problem
095 for both SVMs and RNNs. The Internet represents a vast corpus of knowledge, but the information
096 content is not well-structured for Q&A model training. Recent approaches have wielded question
097 and answer data available from sites such as Yahoo! Answers and Quora [2]. The content from
098 these sites is presented in a format that lends itself to training. Models trained on Q&A site data
099 have similar training objectives, but their space of viable answers is greatly expanded.

100 Yet, the scope of Q&A sites remains small in comparison to the ubiquity of online forums. Forums
101 are structured websites that facilitate question answering, however responses to forum questions
102 need not answer the original question. This leniency makes it difficult to use raw forum data to train a
103 question answering system. Given a forum question, our goal is to distinguish responses that *attempt*
104 *to answer* the question from those that do not. It is a simplification of question-answering in that we
105 are not searching for the best answer to a question. Rather, our work focuses on the general structure
106 of *non-answers* - responses such as comments, follow-up questions, and concluding remarks. By
107 applying our classifier, *non-answers* to forum questions may be removed and the resulting question-
answer pairs used for subsequent model training.

3 Approach

We train a recursive neural network (RNN) that labels responses to forum questions as *answers* or *non-answers*. The *answer* label indicates any response that attempts to answer the question.

Our RNN model is largely based on next-word prediction RNNs. Given the first few words of a sentence, these RNNs predict the next word by assigning a probability distribution over the vocabulary of words. Our model reads each question-response pair from left-to-right and produces a probability distribution over the two classes *answer* and *non-answer*.

3.1 Model

Let x be a question-response pair such that $x = x_1 x_2 \cdots x_n$. For example:

`<Q> What color is the sky? </Q> <R> The sky is blue. </R>`

In the question-response pair above, we have $x_1 = \text{<Q>}$, $x_2 = \text{What}$, $x_3 = \text{color}$, \dots , $x_{13} = \text{</R>}$.

The *answer* and *non-answer* probabilities are given by:

$$P(y = k | x_1 x_2 \cdots x_n) = \text{softmax} \left(U h^{(t)} \right)_k = \hat{y}_k$$

Thus, \hat{y}_1 is the probability of an *answer* and $\hat{y}_2 = 1 - \hat{y}_1$ is the probability of a *non-answer*. We define the recursive function $h^{(t)}$ similar to next-word prediction RNNs:

$$h_t = \text{sigmoid} (H h_{t-1} + L_{x_t})$$

Where L_{x_t} is the column of our word-vector matrix corresponding to the word x_t .

One difficulty with our approach is that we allow for an arbitrary window size. The question-response pair x has n tokens, resulting in n evaluations of h (e.g. t ranges from 1 to n). We assign a label to the question-reponse pair after the final token (e.g. x_n) is seen. Previous work has demonstrated that such models quickly lose memory due to vanishing gradients.

We address this problem in two ways - by limiting t to some fixed value t_{\max} and by extending the model to include Long Short Term Memory as described by Hochreiter, et. al. [3]. In the latter case we rely heavily upon their implementation in Theano.

As noted above, our objective is to maximize the average probability of correctly labeled pairs:

$$\frac{1}{N} \sum_{i=1}^N y^i \cdot \log \hat{y}^i$$

Where N is the number of question-response pairs and y^i is a one-hot vector representing the true label of the i th pair, as recorded by Mechanical Turk workers.

4 Experiment

We trained our RNN classifier using question-response pairs from ten online forums. We examine the classification accuracy as a function of t_{\max} (the maximum number of tokens of the question-response pair taken into account) and iteration number.

We also examine the words that give rise to maximal neuronal activation for the *non-answer* classification label.

4.1 Training Data

Our training data consists of question-response pairs from forum threads. As a simple example, consider the question, "What color is the sky?" and example responses, "The sky is blue.", "I'm not sure.", "It can be either blue or gray.", "Why don't you just Google it?", "I think it's purple." These statements would lead to five training inputs:

162 [1, 0] <Q> What color is the sky? </Q> <R> The sky is blue. </R>
 163 [0, 1] <Q> What color is the sky? </Q> <R> I'm not sure. </R>
 164 [1, 0] <Q> What color is the sky? </Q> <R> It can be either blue or gray. </R>
 165 [0, 1] <Q> What color is the sky? </Q> <R> Why don't you just Google it? </R>
 166 [1, 0] <Q> What color is the sky? </Q> <R> I think it's purple. </R>

167 The first, third and fifth responses are *answers* (label [1,0]) while the second and fourth responses
 168 are *non-answers* (label [0,1]). The beginning and end of each question and answer are delimited
 169 with special tokens. Note that the last response is an incorrect answer, but it still constitutes an
 170 *answer*.
 171

172 4.2 Data Acquisition and Preprocessing

173
 174 We built a forum crawler and downloaded 59,205 question-response pairs from ten online forums.
 175 We used Amazon Mechanical Turk to hand-label each response as an *answer* or a *non-answer*. Turk
 176 workers were instructed to select responses that *attempted to answer* the question. We had two
 177 Turk workers label each response, and kept only those responses in which both workers agreed.
 178 We equalized the number of positively-labeled and negatively-labeled samples resulting in 27,682
 179 question-answer pairs. We randomized the data and then split it into train (80%), dev (10%), and
 180 test (10%) sets. The data is summarized in the table below.

	Percentage of total data set	Number of question-answer pairs
Train	80	22,146
Dev	10	2,768
Test	10	2,768

181
 182
 183
 184
 185
 186 The data is provided as a supplemental attachment. We preprocessed the question-response pairs by
 187 applying the following normalizations:
 188

- 189 • Lower-case the text
- 190 • Remove all non-ASCII characters and punctuation
- 191 • Convert whitespace blocks to single spaces
- 192 • Stem words
- 193 • Remove common (e.g. stop) words
- 194 • Remove words occurring in less than 100 training samples

197 4.3 Evaluation Metric

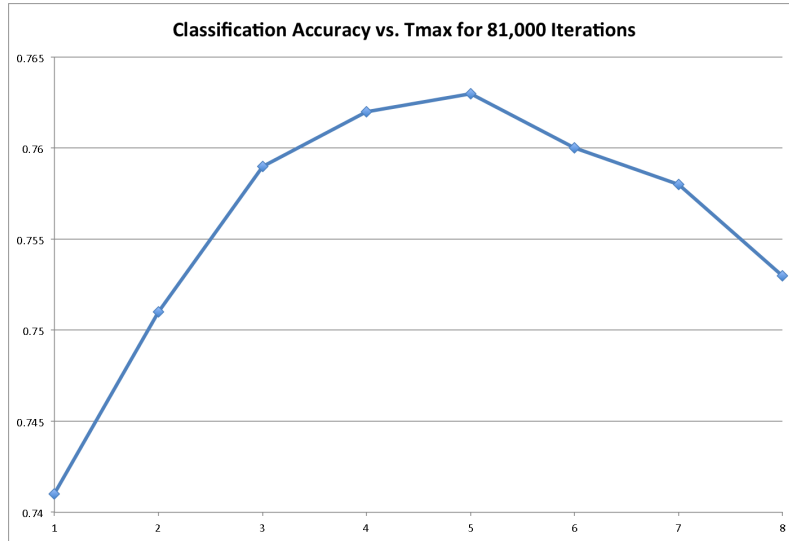
198
 199 We use the number of correctly classified question-response pairs from the test set as our evalua-
 200 tion metric. We also record the words that gave rise to maximal neural activation as a qualitative
 201 evaluation.
 202

203 4.4 Results Without Memory

204
 205 We initially included both the question (<Q> . . . </Q>) and the response (<R> . . . </R>), but early
 206 experiments showed that the response alone was sufficient to classify it as a *non-answer*. We hypothe-
 207 size reasons for this in the discussion. We also found that our model lost memory after $t_{\max} = 5$
 208 recursive steps (Figure 1). Thus we train on the first five words of each response (not counting
 209 the initial <R>). The results below are for our model *without* memory (i.e. no Long Term Short
 210 Term Memory). As shown in Figure 2, the model quickly learned to classify roughly 70% of the
 211 training samples during the first 1,000 iterations. Additional iterations saw gradual performance
 212 improvement with a plateau around 76.3%.

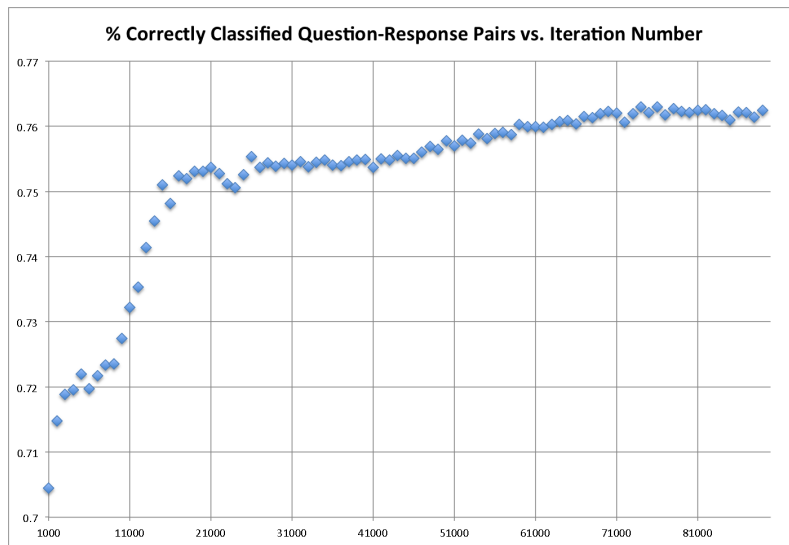
213 In order to qualitatively evaluate the trained model, we fed each vocabulary word into the network
 214 as a one-word sample. Encouragingly, words that were strong *non-answer* features for the SVM
 215 and LR models proved to be strong signals of *non-answers* in the RNN. In particular, the word
 “thanks” produced the highest *non-answer* signal, consistent with the fact that responses starting

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233



234 Figure 1: T_{\max} denotes the number of question-response tokens taken into consideration. Due to
235 gradient vanishing effects, we lose accuracy with $T_{\max} > 5$.
236

237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253



254 Figure 2: Accuracy as a function of training iteration for $T_{\max} = 5$. The model learned quickly
255 during the first thousand iterations, and then slowed.
256

257
258 with “thanks” are very rarely answers. A selection of *non-answer* words that qualitatively validate
259 the model are provided in the table below.
260

Single-Word Input	Non-Answer Probability
thanks	0.83
bump	0.80
nevermind	0.77
inappropriate	0.65

261
262
263
264
265
266
267 Table 1: Words with high non-answer probability. Note that the first three of these single-answer
268 responses are typically written by the original poster. In our SVM and LR models, we include an
269 indicator feature that records whether or not the response was by the original poster.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

4.5 Results with Memory

We attempted to include Long Short Term Memory using the Theano model by Hochreiter, et. al. In theory, this model should have significantly improved our results because it would resolve issues with vanishing gradients, allowing the model to read the entire question-response pair. We started with their sentiment code as a base and attempted to modify it to train our classifier. Unfortunately, we were not able to produce a converging model. The first few iterations produced a minor improvement (67% classification accuracy), but then stagnated at that value. We suspected a bug in our implementation, but ran out of time before we could resolve it.

5 Conclusion

We previously trained a logistic regression classifier and a support vector classifier to evaluate the feasibility of this problem. In those models, we were able to hand-select features beyond the textual data. For example, we indicated the author of the response, the index of the response, and whether the response was quoted by another response. The logistic classifier performed the best with a test-set accuracy of 89.0% using 7-gram features along with several novel features.

In SVM and LR models, we found that the question contributed little to the training accuracy. This result was confirmed by the RNN - we were able to achieve the same RNN classification accuracy using just the response text (`<R> . . . </R>`). This suggests a general structure of *non-answers* to forum questions. Our qualitative evaluation revealed that words such as, “thanks”, “bump” and “nevermind” are strongly indicative of *non-answers*. Note, however, that these words are often written by the original poster, so the post author feature may be an even stronger feature available to the SVM and LR classifiers.

The utility of a response classifier may be demonstrated using a practical example taken from the Apple Technical Support Discussion Forum:

My Apple TV will not connect to WiFi, having put in the pass code it displays connection error -3905. Can anyone help?

1. I too have an ATV2...I could only get it to work properly by uninstalling my McAfee security suite.

RNN Score: 0.21, LR Score: 0.615

2. I get the exact same problem and error code with my NetGear DG384T wireless router; when I try to configure the Wireless network the Apple TV actually lists my wireless SSID but fails to connect when I select it from the list.

RNN Score, 0.10, LR Score: 0.291

3. Take a chill pill, pal! Turn off MAC Address Filtering and you will instantly get a connection. I read a post by a housewife in Australia who nailed the problem - good on yer, Sheila!

RNN Score: 0.11, LR Score: 0.482

4. I was having the same problem with my ATV 2nd Gen. Spoke to an Apple tech, determined that the cause of my problem was wireless interference...Apparently a nearby wifi network was using the same channel that mine was.

RNN Score: 0.33, LR Score: 0.726

5. Thanks everyone for your help, I resolved the issue.

RNN Score: 0.052, LR Score: 0.103

The last answer turns out to be the correct answer. Notice that the 1st and 3rd posts are attempts at answering while the 2nd post is a user who is experiencing the same issue. The last post is by the original poster. Both the RNN classifier and the LR classifier successfully reorder these posts 4, 1, 3, 2, 5. However, the RNN does not classify any of these responses as *answers* (probability < 0.5), likely because the first five tokens of each response is not sufficient to identify them as answers. We can see the strong *non-answer* score of the last response, almost entirely due to the fact that it starts with “thanks”.

324 We found that the question (<Q> . . . </Q>) provided little insight into classifying responses as *non-*
325 *answers*. We hypothesize that this is because *non-answers* generally have a very similar structure
326 (e.g. they might begin with a token like “thanks” or “bump”). Since we were not looking for *correct*
327 answers, the meaning behind each question was not as important as the structure of each response.
328 Put simply, a response such as “Just Google it” is a *non-answer* to a great many questions, whereas
329 it is an *answer* to very few.

330 We previously developed a cross-forum search engine that uses only *answers* (probability > 0.5).
331 The search engine is available at:

332
333 <http://www.gotoanswer.com>

334 Note that this search engine relies upon the logistic regression classifier.

335 Our RNN did not perform as well as the SVM and LR classifiers, but there are quite a few modifica-
336 tions that would likely improve its performance. As a next step, we suggest reducing the question-
337 response space to a specific topic such as computer hardware repair. It is likely that *non-answer* and
338 *answer* responses within a particular subject share common structure. We saw this trend with our
339 SVM and LR classifiers - posts within the Apple Discussion Forum that linked to specific support
340 topics were *answers*, while the phrase “no hackintosh questions” was a common *non-answer*.
341

342

343 **References**

344

- 345 [1] David Azari, Eric Horvitz, Susan Dumais, and Eric Brill. Web-based question answering: A
346 decision-making perspective. In *Proceedings of the Nineteenth conference on Uncertainty in*
347 *Artificial Intelligence*, pages 11–19. Morgan Kaufmann Publishers Inc., 2002.
- 348 [2] Benjamin V Hanrahan, Gregorio Convertino, and Les Nelson. Modeling problem difficulty and
349 expertise in stackoverflow. In *Proceedings of the ACM 2012 conference on Computer Supported*
350 *Cooperative Work Companion*, pages 91–94. ACM, 2012.
- 351 [3] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
352 9(8):1735–1780, 1997.
- 353 [4] Mohit Iyyer, Jordan Boyd-Graber, Leonardo Claudino, Richard Socher, and Hal Daumé III. A
354 neural network for factoid question answering over paragraphs. 2014.
- 355 [5] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for
356 word representation.
- 357 [6] Jun Suzuki, Yutaka Sasaki, and Eisaku Maeda. Svm answer selection for open-domain question
358 answering. In *Proceedings of the 19th international conference on Computational linguistics-*
359 *Volume 1*, pages 1–7. Association for Computational Linguistics, 2002.
- 360 [7] Show-Jane Yen, Yu-Chieh Wu, Jie-Chi Yang, Yue-Shi Lee, Chung-Jung Lee, and Jui-Jung Liu.
361 A support vector machine-based context-ranking model for question answering. *Information*
362 *Sciences*, 224:77–87, 2013.

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377