

Machine Learning Optimizations for SQuAD

Julian Villalpando

Department of Computer Science

Stanford University

Stanford, CA 94305

jvilla32@stanford.edu

Abstract

The baseline codebase offered by the CS224 instructing team offers a solid foundation to build upon. By adding simple, yet powerful upgrades, I was able to improve upon the baseline model in meaningful ways. In this paper, I will detail the implementation and impact of my added features.

1 Introduction

Information retrieval has come a long way – especially in the last few years. Where there used to be teams of 100 hand-crafting feature sets, there are now end-to-end models producing stronger, more robust results. With the recent introduction of attention mechanisms, performances have reached all-time highs. Below, I will detail the techniques that I used to improve upon the baseline model – which include: bidirectional attention flow, feature engineering, and improved span selection

2 Dataset

The Stanford Question Answering Dataset (SQuAD) was generated from over 500 Wikipedia articles and contains over 100,000 question, context, span, and answer pairs. The context represents the body of text, from which the question can be answered. The question has a well defined answer, with the span being the start and ending indices of the answer.

Before developing, I generated several histograms that gave me insight into the general trends of the SQuAD dataset. The most telling feature of the dataset that I found was that most of the answers were relatively short:

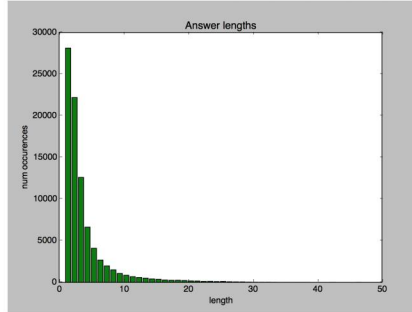


Figure 1: SQuAD Answer lengths

35
36
37
38
39
40
41
42
43
44

It seemed that ~90% of the answers spanned less than 15 indices, which corroborates well with the DrQA paper [1], who capped their answers at 15 in length.

The next thing I noted was that the context lengths hardly, if ever, reached 600 words in size:

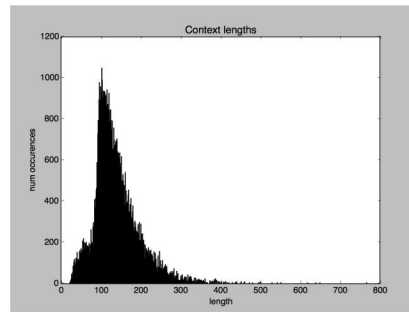


Figure 2: SQuAD Context Lengths

45
46
47
48
49
50
51

Lastly, I measured both the question lengths and the index of the true start position in the context.

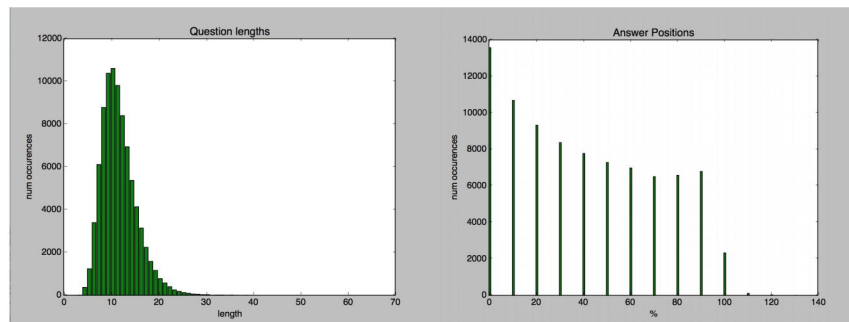


Figure 3: SQuAD Question Lengths and Answer Positions

52
53
54
55

56 **3 Architecture**

57

58 **3.1 Overview**

59

60 My model, like the baseline, is an encoder-decoder neural network. Every
61 word is converted into a vector representation, and then passed into the
62 encoder. This encoder then produces an attention output and distribution,
63 which feeds the decoder. Finally, given the final output distribution, start and
64 end positions are predicted.

65

66 **3.2 Word Representations**

67

68 Every context and question is represented as a matrix of word representations.
69 Every word is converted to a pre-trained, static glove vector, which is then
70 concatenated with a feature vector. The feature vector is quite simple – it
71 represents whether a context word appears in the question.

72

73 **3.3 Encoder**

74

75 The word representations are then passed through a bidirectional recurrent
76 neural net. The forward and backward states are concatenated together, to
77 give the final encodings. The encoder is shared between the context and
78 question encodings.

79

80

81 **3.4 Bidirectional Attention**

82

83 My architecture uses the bidirectional attention mechanism proposed in [2].
84 Given the context and question encodings, the bidirectional attention layer
85 computes a question-to-context (Q2C) and context-to-question (C2Q)
86 attention. It outputs a blended representation of the encoded states as a
87 concatenation of the encoded context, the C2Q output, the encoded context *
88 C2Q output, and the encoded context * Q2C output.

89

90

Equation 1: BiDAF output

91

$$b_i = [c_i; a_i; c_i \circ a_i; c_i \circ c'] \in \mathbb{R}^{8h} \quad \forall i \in \{1, \dots, N\}$$

92

93

94

95

96

97

98

99

100

101

101

3.5 Decoder

The decoder layer takes the blended representations and feeds it through a fully connected layer, followed by a RELU non-linearity. That output is then passed through a downprojecting linear layer, at which point the scores for each context location are softmaxed.

3.6 Loss

102 A cross-entropy loss is taken of both the start and end locations to calculate
103 the loss, and is minimized by the Adam optimizer.

104

105 3.7 Prediction

106

107 At test time, the model predicts the start and end indices – s and e – that
108 maximize: $p_s * p_e$ where $e \leq s + 15$

109

110

111 3.8 Evaluation

112

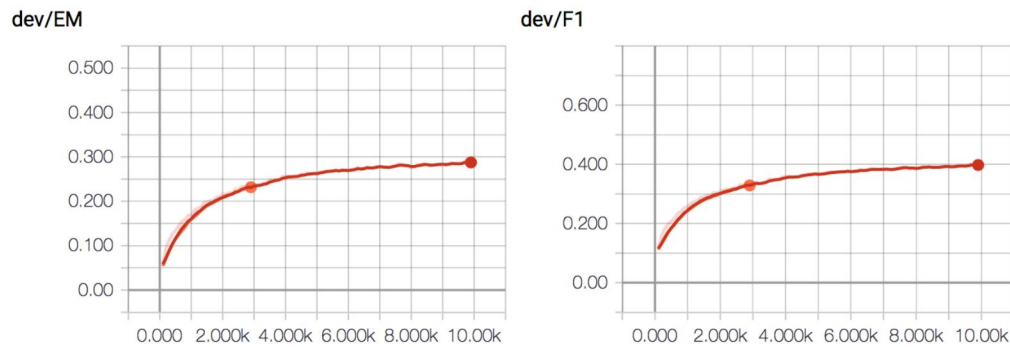
113 F1 and EM scores are used to evaluate the model – where F1 measures
114 precision and recall, while the EM measure exact matching.

115

116 3.9 Hyper parameters

117

118 Context lengths were capped to 400, which showed no significant difference
119 in performance:



120

121

Figure 4: Context lengths 400 vs 600

122

123 Batch sizes were reduced to 32 to accommodate the Microsoft Azure memory
124 limits, which the bidirectional attention flow tended to surpass.

125

126 All else was held constant.

127

128 3.10 Results

129

130 The final results for the model were:

131

132

Table 1: Results

133

	F1	EM
Train	64.472	54.938
Dev	67.643	56.197
Test	67.632	56.572

134

135
136
137
138
139
140
141
142

4 Improvements and Analysis

4.2 Featurization

The simplistic addition of featuring whether or not a context word appears in the question made a significant impact in performance:

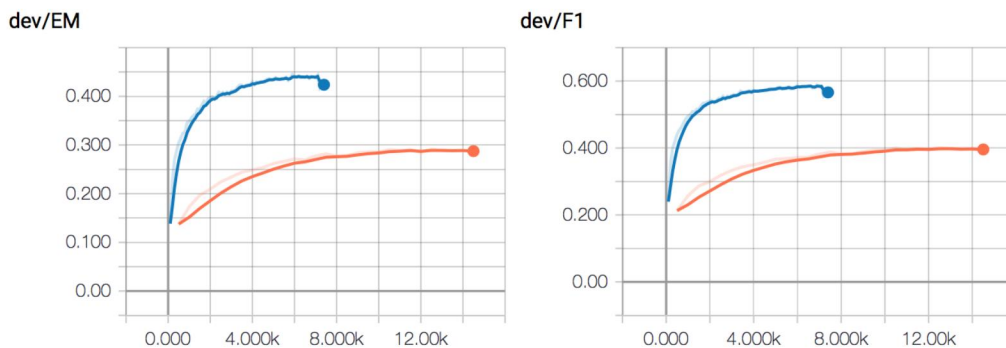


Figure 5: Baseline vs Featurization

143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164

This makes intuitive sense; many of the words in the question are in close vicinity to the answer span in a context question, so the simple feature makes a pronounced difference.

I think that end-to-end models are significantly more powerful than feature-engineered models. We've seen this in class, with regards to the models built by big teams that tried to handcraft all the rules of the English language. End-to-end models have proven not only more powerful, but also significantly easier to develop. However, as seen here, Supplementing features on top of an end-to-end model can give the best of both worlds, and a simple feature can make a pronounced difference.

4.3 Bidirectional Attention

I anticipated that bidirectional attention flow would be my largest performance improvement. That didn't turn out to be the case:

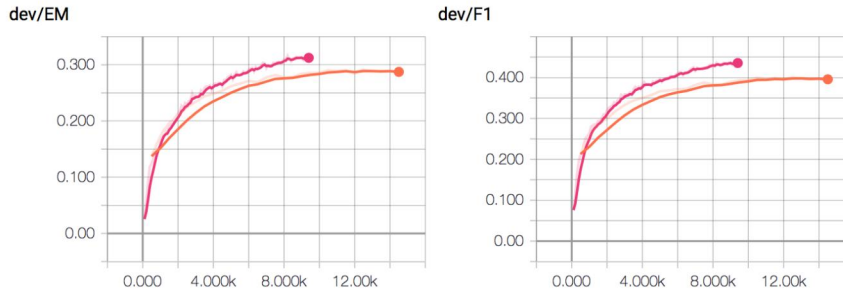


Figure 6: Baseline vs Bidirectional Attention Flow

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180

Upon consulting with peers, I was told that adding a LSTM layer at the end of attention layer, as in the paper, significantly boosted performance. I find this surprising, but understandable. I suppose the attention outputs given in the paper were optimized with the final layer in mind. Thus, their model was able to make better sense of my blended representation detailed in Equation 1.

4.4 Test-time Prediction

One of the biggest problems I noticed in the baseline model, when generating examples, was that the end index would appear before the start index. Given how frequently I saw it, I had guessed that constraining: $start < end < start + 15$ would be a significant improvement. I got the following results:

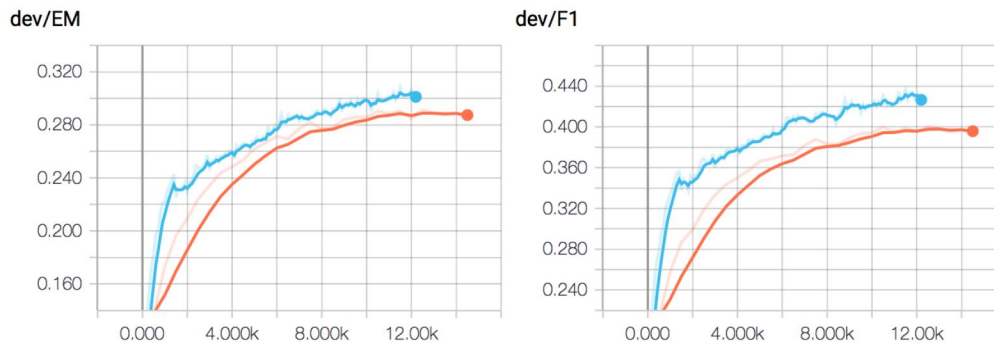


Figure 7: Baseline vs Smarter Span Selection

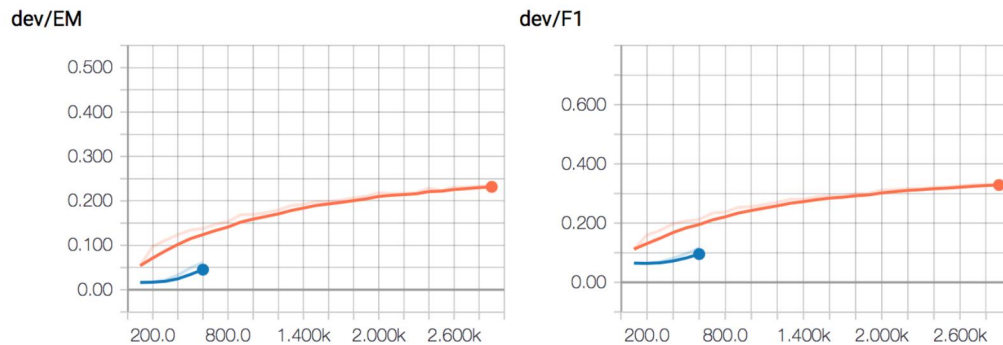
181
182
183
184
185
186
187
188
189
190
191
192
193
194

While the performance improved consistently over all iterations, the improvement was not as large as I thought it would be. I suppose might occur because the start and end distributions are still calculated separately, and the loss + gradient is determined from them. Meanwhile, the F1 and EM scores, are calculated from the joint probability, so the weights were not directly optimizing for F1 and EM. Nevertheless, the improvements were significant and thus kept for the final implementation. I think conditioning the end distribution on the start distribution – as in the Answer Pointer implementation [4] - would have paired nicely with this heuristic, but I was unable to successfully implement it as described in the handout.

195 4.5 Self Attention

196

197 In my attempt to improve the model, I implemented self-attention as proposed
198 in [3]. However, given the memory constraints, I had to reduce batch sizes to
199 20, context lengths to 350, and hidden sizes to 100 to run the model on
200 Microsoft Azure. I believe that given these constraints, I was unable to
201 produce a self-attention model that surpassed the less memory-intensive
202 baseline
203



204

205

Figure 8: Baseline vs Self Attention

206

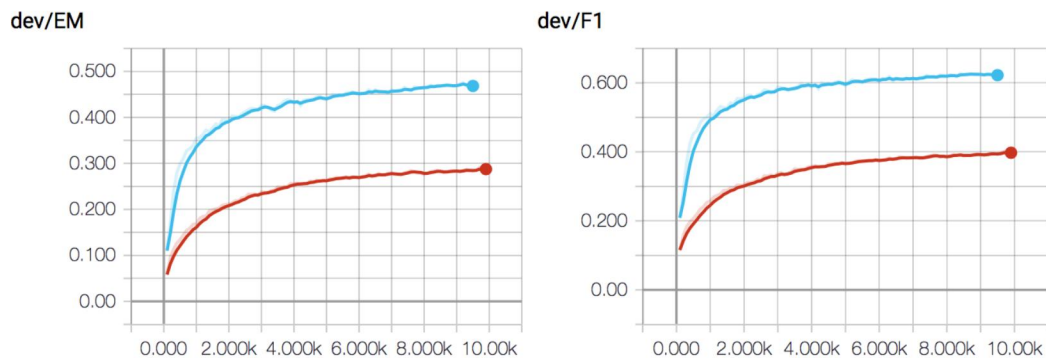
207

208 4.6 Final Product

209

210 As detailed in the Table 1, my final performance was:

211



212

213

Figure 9: Baseline vs Final Model

214

215 My features seemed to pair well on top of each other. The performance of the
216 final model suggests that the improvements all added on to each other, with
217 little overlap – as in output final = improvement from A + improvement from
218 B + improvement from C. I would wager that the improvement of the start
219 and end predictions did not overlap with the improvements of the word
220 embedding and the attention mechanism, and that most of the overlap
221 occurred between the latter two. This is because both the word embedding and
222 attention mechanism affect the final blended representation of the hidden and
223 context states.

224

225 **5 Conclusion**

226

227 Neural nets have pushed the limits of question answering – with several
228 leading models now surpassing human answering in F1 scoring. That being
229 said there is still a long way to go. For my model in particular, I'd like to
230 spend more time making the start and end positions dependent upon each
231 other – as alluded to previously. I think I'd start with a simpler model than the
232 Answer Pointer described in [4], to first verify that it would pair well with my
233 current model, and then proceed to implement more advanced iterations of it.

234 **References**

- 235 [1] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer
236 open-domain questions. arXiv preprint arXiv:1704.00051, 2017.
- 237 [2] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional
238 attention flow for machine comprehension. arXiv preprint arXiv:1611.01603, 2016.
- 239 [3] Natural Language Computing Group, Microsoft Research Asia. “R-NET: Machine: Reading
240 Comprehension With Self-Matching Networks.” Microsoft.
- 241 [4] Shuohang Wang and Jing Jiang. Machine comprehension using match-1st mand answer pointer.
242 arXiv preprint arXiv:1608.07905, 2016.