

Machine Comprehension for SqUAD dataset

Vikas Bahirwani
Microsoft Corporation
vikasb@stanford.edu

Erika Menezes
Microsoft Corporation
emenezes@stanford.ed

Abstract

We focus on predicting the start and end indices of the answers. Our approach explores the effectiveness of RNNs, BiRNNs, LSTMs, BiLSTMs etc. We explore the use of Attention in addressing this problem as well. We also observe that our model is underfitting and our next steps would have been to develop a more complex model to overcome the same.

1 Introduction

Machine comprehension using Deep Learning models is a growing field in Natural Language Processing. While it possesses immense potential it also presents a lot of challenges - special the challenge to design a neural network to fit the downstream task at hand.

In this assignment, we address the reading comprehension task of generating answers from context paragraphs given the questions. Our approach starts with a simple baseline model using RNNs to capture question and context knowledge.

We then update our model to use LSTMs - which help us learn longer paragraphs and address the answer generation problem better than RNNs. We focus on predicting the start and end indices of the answers.

2 Background/Related Work

There have been many deep learning models proposed for machine comprehension. Wang et al.[1] work is based on the assumption that a span in a passage is more likely to be the correct answer if the context of this span is very similar to the question. The novelty in this paper is the Multi-Perspective Context Matching (MPCM) model that identifies the answer span by matching the context of each point in the passage with the question from multiple perspectives. There is also the work by Xiong et al.[2] that focuses on how to recover from local maxima from incorrect answers. The Dynamic Co-attention Network combines the question and the document in order to focus on relevant parts of both. Then a dynamic pointing decoder iterates over potential answer spans.

3 Approach

3.1 Data Analysis

To better understand the problem, we visualized the data by printing it from `qa_answer.generate_answers()`. Once the dataset has been read into a list of tuples of (context, question, question_uuid) we used this and the vocabulary passed as an argument to

41 generate the answer for a given a_s and a_e . We are still working on encoding and decoding
42 in order to generate model predictions.

43

44 **3.2 Approach**

45

46 Our first approach was to have a simple encoder decoder model as a baseline. This was
47 implemented using an LSTM over the question, and another LSTM for the context by using
48 the initial state as the final state from the question. We then used the final states from the
49 question and paragraph through a 1 layer neural network to predict the start and end of the
50 answer. The following steps capture the approach:

51

- 52 1. Question \rightarrow LSTM \rightarrow Q
- 53 2. Paragraph \rightarrow LSTM(initial = Q) \rightarrow P
- 54 3. KRep = [Q,P]
- 55 4. $a_s = \text{softmax}(\text{KRep} * W1) + B1$
- 56 5. $a_e = \text{softmax}(\text{KRep} * W1) + B1$

57

58 We realise that the biggest drawback of this model is that it does not include attention and
59 in order to fix this we come up with a slightly more complex model that involves the
60 following steps:

61

- 62 1. Question \rightarrow LSTM \rightarrow Q
- 63 2. Paragraph \rightarrow LSTM \rightarrow P
- 64 3. $A = \text{softmax}(P Q^T)$ //Compute context vector for Q \rightarrow P
- 65 4. $C_P = A Q$ // and mix with P
- 66 5. $P = \text{concat}(C_P, P) W + b$ // Mix it with P (Krep)
- 67 6. $a_s = \text{softmax}(\text{KRep} * W1) + B1$
- 68 7. $a_e = \text{softmax}(\text{KRep} * W1) + B1$

69

70 The last approach that we tried was to introduce non linearity in the decoder by using a
71 ReLU activation function in the neural network .

72

- 73 1. Question \rightarrow LSTM \rightarrow Q
- 74 2. Paragraph \rightarrow LSTM \rightarrow P
- 75 3. $A = \text{softmax}(P Q^T)$ //Compute context vector for Q \rightarrow P
- 76 4. $C_P = A Q$ // and mix with P
- 77 5. $P = \text{concat}(C_P, P) W + b$ // Mix it with P (Krep)
- 78 6. $a_s = \text{softmax}(\text{KRep} * W1) + B1$
- 79 7. $a_e = \text{softmax}(\text{KRep} * W1) + B1$

80

81 **4 Experiments**

82

83 We experimented (on all approaches) both locally and on GPU. Locally we used 1000
84 samples (from training) to train the model across 10 epochs (batch size 10). At each epoch,
85 we calculated F1 on 50 validation dataset samples.

86

87 In addition, we also evaluated the F1 score of the overall model.

88

89 Here are the loss and F1 scores from the final approach.

90

91 Epoch 1 out of 10

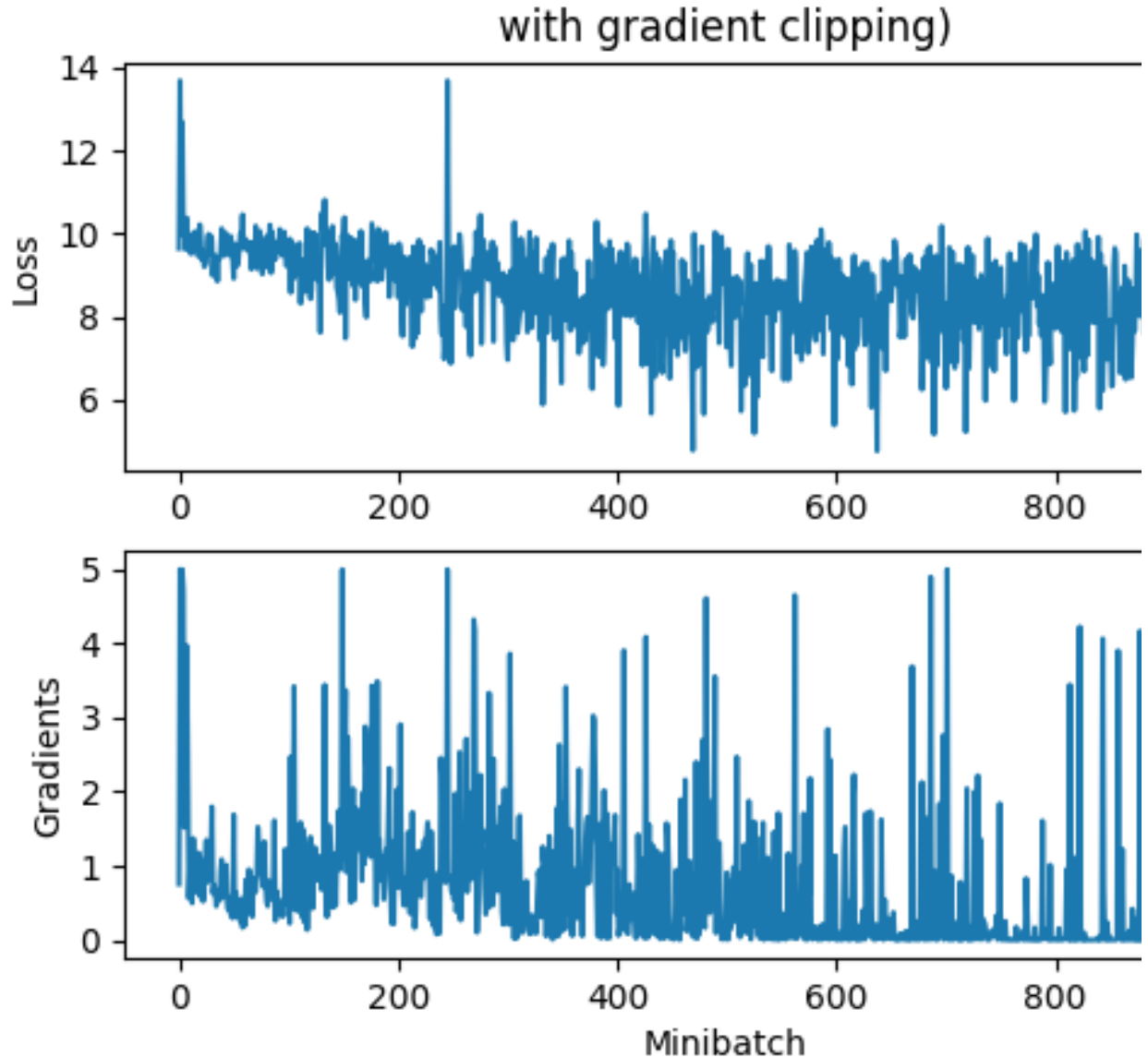
92 train loss: 9.7127

93 Score 6.181432 , best_score so far 6.181432

94

95 Epoch 2 out of 10

96 train loss: 9.3238
97 Score 4.153968 , best_score so far 6.181432
98
99 Epoch 3 out of 10
100 train loss: 8.8529
101 Score 6.735965 , best_score so far 6.735965
102
103 Epoch 4 out of 10
104 train loss: 8.5121
105 Score 4.476740 , best_score so far 6.735965
106
107 Epoch 5 out of 10
108 train loss: 8.3757
109 Score 2.268926 , best_score so far 6.735965
110
111 Epoch 6 out of 10
112 train loss: 8.2868
113 Score 4.464495 , best_score so far 6.735965
114
115 Epoch 7 out of 10
116 train loss: 8.2392
117 Score 3.814750 , best_score so far 6.735965
118
119 Epoch 8 out of 10
120 train loss: 8.2109
121 Score 4.473124 , best_score so far 6.735965
122
123 Epoch 9 out of 10
124 train loss: 8.2412
125 Score 4.029124 , best_score so far 6.735965
126
127 We noticed that we used our models were underfitting because our loss will not go down.
128 (As opposed to overfitting where loss on train is almost 0 or F1 is very high but on
129 validation set the performance degrades).
130
131 Furthermore, to confirm that we were not falling into the trap of gradient explosion, we
132 generated the following plot. Essentially it was the same with and without clipping.
133 (GradNorm 5)
134



135

136

137

5 Conclusion

138

139

140

141

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle that the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

142

References

143

144

145

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D. S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609-616. Cambridge, MA: MIT Press.

146

147

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.

148

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory

149 recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of*
150 *Neuroscience* **15**(7):5249-5262.