
CS 224N Default Final Project: Question Answering

Raghunath Krishnamurthy(kichyrag)
SCPD - Stanford University
kichyrag@stanford.edu

Abstract

The default project is to build a system for the Stanford Question Answering Dataset (SQUAD). The goal is to understand the context from a provided paragraph, followed by query(question) and system should output the correct answer. Basic structure of the system implemented will be improvements to the baseline model, which has **RNN encoder layer**, which encodes both the context and the question into hidden states, an **Attention Layer**, that combines the context and an output layer, which applies a fully connected layer and then two separate softmax layers, one of them to get the start location, and the other one to determine the end location of the answer. The improvements attempted, were to implement **Bidirectional attention flow (BiDAF)**[1], which was quite unsuccessful because of decode layer was kept same softmax, and this impacted the dev F1 and EM score to be very low (2.6%) and failed to improve even after 10k iterations. Finally, several adjustments to the baseline model was done such as dropout, learning rate decay and LSTM instead of GRU. Finally the model with LSTM and dropout of 0.20 achieved the best results. Future work is to improve on the BiDAF and implement LSTM instead of softmax for decoding to achieve higher accuracy.

1 Introduction

The intuition behind attention mechanism is that the model is able to recognize the difference between context encodings, where the question pays more attention.

Some examples of the Questions and the true answers:

```
CONTEXT: (green text is true answer,agenta background is predicted start,red background is predicted end,underscores are unknown tokens). Length: 188
not only are all the major british architects of the last four hundred years represented, but many european ( especially italian ) and american architects ' drawings are
held in the collection . the riba 's holdings of over 330 drawings by andrea palladio are the largest in the world , other europeans well represented are jacques
gentilhatre , and antonio gaudi . british architects whose drawings , and in some cases models of their buildings , in the collection , include : inigo jones , sir
christopher wren , sir john vanbrugh , nicholas hawkmoor , william kent , james gibbs , robert adam , sir william chambers , james wyatt , henry holland , john nash , sir
john soane , sir charles barry , charles robert cockerell , augustus welby northmore pugin , sir george gilbert scott , john loughborough pearson , george edmund street ,
richard norman shaw , alfred waterhouse , sir edwin lutyens , charles rennis mackintosh , charles holden , frank hoar , lord richard rogers , lord norman foster , sir
nicholas grishaw , zaha hadid and alick horsnell .
QUESTION: which architect , famous for designing london 's st. paul cathedral , is represented in the riba collection ?
TRUE ANSWER: sir christopher wren

CONTEXT: (green text is true answer,agenta background is predicted start,red background is predicted end,underscores are unknown tokens). Length: 112
the crew of apollo 8 sent the first live televised pictures of the earth and the moon back to earth , and read from the creation story in the book of genesis , on christmas
eve , 1968 , an estimated one-third of the population of the world saw-either live or delayed-the christmas eve transmission during the ninth orbit of the moon . the
mission and christmas provided an inspiring end to 1968 , which had been a troubled year for the us , marked by vietnam war protests , race riots , and the assassinations of
civil rights leader martin luther king , jr. , and senator robert f. kennedy .
QUESTION: how much of the population of earth ended up seeing the images of the earth and the moon ?
TRUE ANSWER: one-quarter

CONTEXT: (green text is true answer,agenta background is predicted start,red background is predicted end,underscores are unknown tokens). Length: 38
in early 2012 , nfl commissioner roger gooden stated that the league planned to make the 50th super bowl " special " and that it would be " an important game for us as a
league " .
QUESTION: what one word did the nfl commissioner use to describe what super bowl 50 was intended to be ?
TRUE ANSWER: spectacular

CONTEXT: (green text is true answer,agenta background is predicted start,red background is predicted end,underscores are unknown tokens). Length: 130
to remedy the causes of the fire , changes were made in the block ii spacecraft and operational procedures , the most important of which were use of a nitrogen/oxygen
mixture instead of pure oxygen before and during launch , and removal of flammable cabin and space suit materials . the block ii design already called for replacement of the
block i plug-type hatch cover with a quick-release , outward opening door . nasa discontinued the manned block i program , using the block i spacecraft only for unmanned
saturn v flights . crew members would also exclusively wear modified , fire-resistant block ii space suits , and would be designated by the block ii titles , regardless of
whether a lm was present on the flight or not .
QUESTION: what type of materials inside the cabin were removed to help prevent more fire hazards in the future ?
TRUE ANSWER: flammable cabin and space suit materials
```

2 Approach

Approach 1

Basic structure of the system implemented will be improvements to the baseline model, which has **RNN encoder layer**, which encodes both the context and the question into hidden states, an **Attention Layer**, that combines the context and an output layer, which applies a fully connected layer and then two separate softmax layers, one of them to get the start location, and the other one to determine the end location of the answer.

The context and question are fed into an encoder, which had two independent BiLSTM with dropout. The context is represented as sequence of d-dimensional word embedding, and the question by a sequence of d-dimensional word embedding. The GloVe embedding are fed into a 1-layer bidirectional LSTM, which produces and sequence of hidden states, both forward and backward.

Attention layer is a dot-product attention, with context hidden states attend to question hidden states, which produces attention distribution. The attention distribution is used to produce weighted sum of question hidden states.

The attention output is then concatenated to context hidden states to obtain blended representation.

Decoder The fully connected layer is then fed into two separate softmax layers, one of them to get the start location, and the other one to determine the end location of the answer.

Approach 2

The next model attempted to improve on the basic attention, instead of using dot-product is using BiDAF [1], which uses similarity matrix to compute context-to-question attention using weighted sum of question hidden states with softmax on similarity matrix. The next step is perform question-to-context attention, which takes max of the corresponding row of the similarity matrix. The attention distribution here would be softmax of max of similarity matrix, which is used to take a weighted sum of context hidden states.

Finally the attention output is determined by combining context hidden states, context-to-question attention output and question to context attention output. This is then fed to the decoder from approach 1.

3 Experiments

Based on the readings and recommendation, BiDAF would have been the best model. The model suffered from decoder as it was fed into the softmax and needed something like Bidirectional LSTM to have had good accuracy.

I decided to change for tuning better hyperparameters to fit the basic attention model to perform better. Couple of the experiments that were undertaken were 1. Change the Batch size (Increase the batch size)

2. Increase the size of hidden states 3. increase the embedding size of the pre-trained word vectors

Most of the above changes worked well on training, but did not work well on official evaluation and suffered OOM on GPU.

Finally the best combination found by changing the dropout to 0.20, which performed well on dev and test set. I was able to achieve 44% F1 accuracy with that setting.

4 Conclusion and Next steps

The basic dot-product attention structure works well in achieving good results, and the key takeaway was just increasing the embedding size, batch size or hidden states is not good. The important future works are definitely adding regularization to prevent over-fitting. I am planning on continuing to improve and work on completing the BiDAF decoder and other models like co-attention to improve the accuracy.

References

[1] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional attention flow for machine comprehension. arXiv preprint arXiv:1611.01603, 2016.