



QUESTION ANSWERING

EXAMINING THE TOP TRENDS AND METHODS FROM THE SQUAD TASK

INTRODUCTION

- Extractive question answering is a machine reading comprehension problem which has exploded in popularity with the release of the Stanford Question Answering Dataset (SQuAD)
- Bi-Directional Attention Flow for Machine Comprehension (BiDAF) is a well known and good performing architecture for this task
- Self Attention is a promising component used in a different popular architecture R-Net
- We aim to combine these to build a high performing architecture for the SQuAD task. We investigate the importance of different components and experiment with transfer learning using ELMo contextual embeddings.

PROBLEM

- Given a context paragraph C, and a question Q, provide the token span C[i:j] which answers the question.
- There are two metrics reported for this task:
 - Exact Match (EM) measures the percentage of questions that received the exact correct answer span prediction
 - F1 gives a averaged measure of overlap between the predicted and ground truth answer spans. To calculate this both are treated as a bag of tokens.

DATASET

- SQuAD v1.1 consist of 107,785 question-answer pairs on 536 articles from Wikipedia where the answer to every question is a span from the associated article passage

Example: Similarly, it is not known if **L** (the set of all problems that can be solved in logarithmic space) is strictly contained in P or equal to P. Again, there are **many complexity classes** between the two, such as NL and NC, and it is not known if they are distinct or equal classes.

What lies between L and P that prevents a definitive determination of the relationship between L and P?
 Ground Truth Answers: complexity classes, many complexity classes, many complexity classes

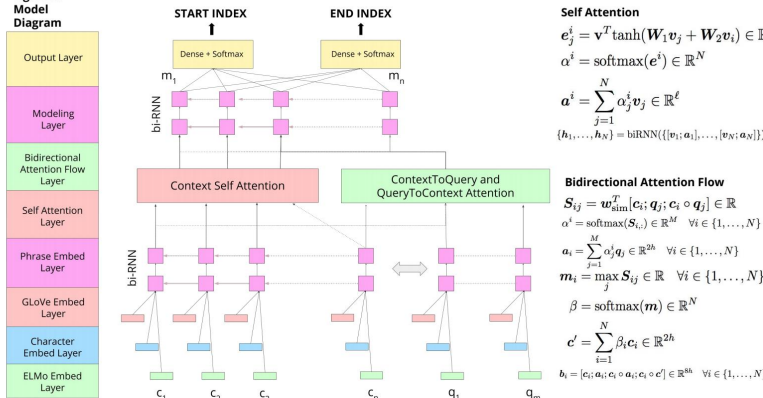
What variable is associated with all problems solved within logarithmic space?
 Ground Truth Answers: L, L, L

Figure 1. Question Breakdown



MODEL ARCHITECTURE

Figure 2. Model Diagram



Self Attention

$$e_j^i = v^T \tanh(W_1 v_j + W_2 v_i) \in \mathbb{R}$$

$$\alpha^i = \text{softmax}(e^i) \in \mathbb{R}^N$$

$$a^i = \sum_{j=1}^N \alpha_j^i v_j \in \mathbb{R}^d$$

$$(h_1, \dots, h_N) = \text{biRNN}([v_1; a_1], \dots, [v_N; a_N])$$

Bidirectional Attention Flow

$$S_{ij} = w_{\text{sim}}^T [c_i; q_j; c_i \circ q_j] \in \mathbb{R}$$

$$\alpha^i = \text{softmax}(S_{i,:}) \in \mathbb{R}^M \quad \forall i \in \{1, \dots, N\}$$

$$a_i = \sum_{j=1}^M \alpha_j^i q_j \in \mathbb{R}^{2k} \quad \forall i \in \{1, \dots, N\}$$

$$m_i = \max_j S_{ij} \in \mathbb{R} \quad \forall i \in \{1, \dots, N\}$$

$$\beta = \text{softmax}(m) \in \mathbb{R}^N$$

$$c' = \sum_{i=1}^N \beta_i c_i \in \mathbb{R}^{2k}$$

$$b_i = [c_i; a_i \circ a_i; c_i \circ a_i] \in \mathbb{R}^{3k} \quad \forall i \in \{1, \dots, N\}$$

RESULTS

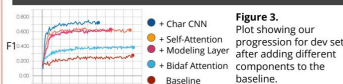


Figure 3. Plot showing our progression for dev set after adding different components to the baseline.

Similarly, it is not known if **L** (the set of all problems that can be solved in logarithmic space) is strictly contained in P or equal to P. Again, there are **many complexity classes** between the two, such as NL and NC, and it is not known if they are distinct or equal classes.

What lies between L and P that prevents a definitive determination of the relationship between L and P?
 True Answer: **complexity classes**. Predicted Answer: **complexity classes**

In 2007, the Kenyan government unveiled vision 2030, an economic development programme it hopes will put the country in the same league as the Asian economic tigers by the year 2030 ...

What did Kenya reveal in 2030?
 True answer: **vision 2030**. Predicted answer: **asian economic tigers**



Figure 4. Model metrics(dev) on different question types.



Figure 5. Model metrics(dev) by context length groups.

CHAR CNN

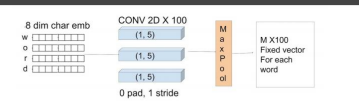


Figure 6. Char CNN to learn high dimensional char level word embeddings

ABLATION

	F1	EM
Baseline	.26	.19
Baseline ++	.31	.21
Full(batch=100)	.78	.63

	F1 (batch=64)	EM
Full + ELMo	.71	.57
Self Attention	.71	0.55
Modeling Layer	.44	.32
Character CNN	.65	.52
BiDAF	.53	.38

Figure 7. Performance on dev (batch size 100)

Figure 8. Ablation Study with smaller batch size (64) than our final model on dev set to understand how removing certain components affected the final model.

- The modeling layer provides a significant performance boost
- We need to fix our ELMo integration since it is negatively affecting performance. It may be better to pass directly to the modeling layer
- Self attention adds only marginal benefit. Most likely due to the fact that the context only attends to itself.
- Character embeddings were more beneficial than we expected
- BiDAF was more beneficial for this task

ANALYSIS

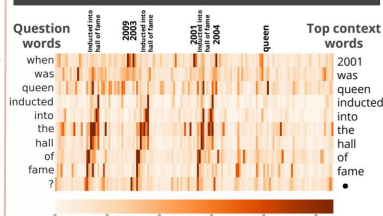


Figure 9. Bidirectional attention visualization using similarity matrix between question and context.

- Key takeaways of hyperparameter tuning experiments**
- We found that the dropout rate of 0.2 gives performance boost but overfits as seen in figure 7 and 8
 - Model learns slower with learning rate of 0.0001 as seen in Figure 9
 - Batch size of 100 gives better performance than 32 and 64

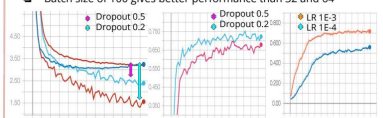


Figure 10. Loss vs iterations for different dropouts (dev). Figure 11. F1 vs iterations score for different dropouts (dev). Figure 12. F1 vs iterations for different learning rates (dev).

CONCLUSION

- The modeling layer is very important to the performance of this model
- The model needs more sophisticated language understanding to handle SQuAD 2.0 and abstractive answering.
- Character embeddings helped more than expected and make the model robust for out of vocabulary tokens.
- Attention is very powerful and efficient. RNNs are slow to train.

FUTURE WORK

- Output layer that conditions the end prediction on the start prediction
- Highway Layers that are used in the BiDAF model
- Experimenting more in how to integrate ELMo and new transfer learning methods such as BERT
- Inputting more features (POS tags, NE tags, EM tags, TRIFD, etc.)
- Test different optimization methods
- Use non-RNN approach (CNNs, Transformer)

REFERENCES

- Mijoon Seo, Anandthi Kambhaji, Ali Farhadi, and Hanmaneh Hajishirzi. Bidirectional attention flow for machine comprehension. arXiv preprint arXiv:1611.01603, 2016.
- Natural Language Computing Group, Microsoft Research Asia, R-NET: MACHINE READING COMPREHENSION WITH SELF-MATCHING NETWORKS
- CS224N Final Project Worksheet