

Question-Answering

Stanford Attentive Reader

Encoder:

Bi-LSTM

Apply encoding to context

$P = \{p_1, p_2, \dots, p_n\}$ through a bi-LSTM network.

$\tilde{h} \rightarrow$ hidden size.

$$\vec{h}_i = \text{LSTM}(\vec{h}_{i-1}, p_i), i=1, \dots, n$$

$$\overleftarrow{h}_i = \text{LSTM}(\overleftarrow{h}_{i+1}, p_i), i=n, \dots, 1$$

and $\tilde{p}_i = \text{concat}(\vec{h}_i, \overleftarrow{h}_i) \in \mathbb{R}^h$

where $h = 2\tilde{h}$.

The resulting forward and o/p
vectors $\overrightarrow{H_p} = \{(\overrightarrow{h_p})_i\}$ and

$$\overleftarrow{H_p} = \{(\overleftarrow{h_p})_i\} \text{ are}$$

concatenated into a final
paragraph representation

$$(H_p)_i = (\overrightarrow{(H_p)_i}, \overleftarrow{(H_p)_i}).$$

Each o/p vector is chosen to be in

$$R^{\tilde{h}} \quad (\underline{\tilde{h} = 200}). \text{ Thus,}$$

$$\text{concatenation is in } \underline{R^{2\tilde{h}}(R^h)}$$

Another (bi-LSTM or bi-GRU)

is applied to map the

question sequence $q = \{q_1, q_2, \dots, q_m\}$

into $h_q = (\overrightarrow{h_q}_m, \overleftarrow{h_q}_m)$

~~* GRU~~ GRU is computationally cheaper than LSTM.

All of the o/p states are used in paragraph representation but only the final states are used in question representation.

Attention:

In this step, goal is to ~~create~~ compare the question embedding and all the contextual (para) embeddings, and select the pieces of information that are relevant to the question. Computation of a probability distribution α depending on the degree of relevance b/w word p_i and question q and then produce an op vector \odot which is a weighted combination of all contextual embeddings $\{P_i\}$:

$$\alpha_i = \text{softmax}_i q^T W_s \tilde{p}_i = \text{softmax}(q^T W_s (H_{\tilde{p}_i}))$$

$$o = \sum_i \alpha_i \tilde{p}_i$$

$i=1, \dots, n$

$W_s \in \mathbb{R}^{h \times h}$ is a bilinear term, which allows to compute a similarity b/w q and \tilde{p}_i more flexibly than with just dot product:

Prediction:

Using the off vector o , the system outputs the most likely answer using:

$$a = \arg \max_{a \in \mathcal{A}} W_a^T o$$

finally, the system adds a softmax
fn on top of $W a^T$ and adopts
a negative log-likelihood objective
for training.