Graph-Structured Representations for Visual Question Answering

Damien Teney, Lingqiao Liu, Anton van den Hengel

Overview



Network architecture



Technical details

Propagation of information over the graph from neighbours, over several iterations

$$h_i^0 = 0$$

$$n_i = \text{pool}_j (e'_{ij} \circ x'_j)$$

$$h_i^t = \text{GRU} (h_i^{t-1}, [x'_i; n_i]) \qquad t \in [1, T].$$

Matching the graphs of questions and image

Attention weights

$$a_{ij} = \sigma \left(W_5 \left(\frac{x_i^{'\mathsf{Q}}}{\|x_i^{'\mathsf{Q}}\|} \circ \frac{x_j^{'\mathsf{S}}}{\|x_j^{'\mathsf{S}}\|} \right) + b_5 \right)$$

- Weighted sum of the features
$$y_{ij} = a_{ij} \cdot [x_i^{'' Q}; x_j^{'' S}]$$

Results



Is the woman exercising '









Does he walk like an idiot ? Answer: yes







Answer: no yes





Who is sitting between toys ? Answer: baby





Answer: brown

As seen on P/R curve: the model's output (after softmax) or sigmoid) is a good measure of its confidence/uncertainty, especially when trained with soft scores as targets.

Practically, this can be used to derive the answer I don't know.



THE UNIVERSITY ofADELAIDE

human human bee doud sun monkee





Is the man sitting on the armrest?



What is underneath the arched Answer: rug plant



What color are the pillows on the

