

2.2 Attention

Let's discuss attention in terms of machine translation. Some problems with the standard encoder-decoder architecture are that long-term dependencies between words that are hard to learn, and that the hidden state attempts to store information of any length into a hidden vector of fixed size. In fact, when generating the next translated word in a sentence, we do not necessarily need to consider the entire input. We can instead choose to only attend to relevant words in an input.

In a vanilla sequence to sequence model, we generate a next word $y^{(t)}$ using the hidden state $s^{(t)}$, which is a function of $s^{(t-1)}$ and $y^{(t-1)}$. With attention, however, we consider $s^{(t)}$ as a function of three things: $s^{(t-1)}$, $y^{(t-1)}$, and some *context* $c^{(t)}$. The context vector is a weighted sum of all input activations $a^{(t')}$:

$$c^{(t)} = \sum_{t'} \alpha^{(t,t')} a^{(t')} \quad (20)$$

where $\alpha^{(t,t')}$ is the amount of attention that $y^{(t)}$ should pay to $a^{(t')}$. Omitting the softmax for simplicity, one way to compute attention is to train a small neural network to compute $\alpha^{(t,t')}$ given $s^{(t-1)}$ and $a^{(t')}$.

Another way of thinking of attention is that given a *query* \mathbf{q} and a set of *key, value* pairs (\mathbf{k}, \mathbf{v} pairs), we compute a context as a weighted sum of the values, where weights are computed using compatibility function between the query and each key. In the scenario described above, the query is the last decoder state $s^{(t-1)}$, and the keys and values are the same and equal the encoder states (activations) $a^{(t')}$. As for computing the compatibility function between \mathbf{q} and \mathbf{k} , there are several ways:

- **Dot product.** $\text{attn_score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$. This method forces input and output encodings to be in the same space.
- **Scaled dot product.** Scale of dot product increases with larger dimensions, so $\text{attn_score}(\mathbf{q}, \mathbf{k}) = \frac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$
- **Bilinear function.** Relaxing the requirement for input and output encodings to be in the same space, we have $\text{attn_score}(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W_a \mathbf{k}$.
- **Multi-layer perceptron (additive similarity).** $\text{attn_score}(\mathbf{q}, \mathbf{k}) = \mathbf{w}_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$.