# Attention Revisited

## Seq 2 Seq with Attention



$h_1 \; h_2 \; - \; h_N$

attention scores $e^t$ :  $\quad e^t = [s_t^\top h_1 \; - \; . \; s_t^\top h_N] \in \mathbb{R}^N$

$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$  take softmax
to get the attention distribution $\alpha^t$
for this step.

$$a_t = \sum_{i=1}^{N} \alpha_i^t \, h_i \in \mathbb{R}^h$$

concatenate  $[a_t ; s_t] \in \mathbb{R}^{2h}$

and proceed as in the non-attention
seq2seq model

Attention: to focus on certain parts of the source

solves the ==bottleneck problem==

helps with the vanishing gradient problem

provides some interpretability

⤵

==get (soft) alignment for free==

There are several attention variants

values  $h_1, h_2, \cdots, h_N \in \mathbb{R}^{d_1}$    <span style="color:blue">similar to (encoder hidden state)</span>

a query    $S \in \mathbb{R}^{d_2}$    <span style="color:blue">(similar to decoder hidden state)</span>

Attention
{
① ==computing attention scores==

   ==$e \in \mathbb{R}^N$==

② Taking softmax to get attention distribution $\alpha_i$

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

③ Using attention distribution to take weighted sum of values:

$$a = \sum_{i=1}^{N} \alpha_i h_i \in \mathbb{R}^{d_1}$$
}

$\implies$ attention output $a$

(sometimes called the context vector)

Several ways to compute $e \in \mathbb{R}^N$ from $h_1 \cdots h_N \in \mathbb{R}^{d_1}$

$s \in \mathbb{R}^{d_2}$

① Basic dot product attention

$$e_i = s^T h_i \in \mathbb{R}$$

assumes $d_1 = d_2$

② multiplicative attention:

$$e_i = s^T W h_i \in \mathbb{R}$$

Where $W \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix

③ Reduced rank multiplicative attention:

$$e_i = s^T (U^T V) h_i$$

$$= (Us)^T (V h_i)$$

low rank matrices $U \in \mathbb{R}^{k \times d_2}$, $V \in \mathbb{R}^{k \times d_1}$

$$k << d_1, d_2$$

④ Additive attention: $e_i = V^T \tanh(W_1 h_i + W_2 s)$

$\in \mathbb{R}^{d_3 \times d_1}$

$\in \mathbb{R}^{d_3}$

$\in \mathbb{R}$

$\in \mathbb{R}^{d_3 \times d_2}$

Given a set of vector values,
and a vector query, attention is a technique