

Introduction to Self Attentionwith trainable weights

(aka scaled dot product attention)

- \* Raw dot product b/w embeddings can give incorrect or ambiguous attention weights.

Ex,

The dog chased the ball but couldn't catch it.

→ "it" should attend more to ball.

→ raw dot product may give equal scores.

→ model cannot fix this without training.

- \* Instead of directly comparing input embedding, we do transform first and then compute attention in a learned space. This introduces trainable parameters.

Transforming embedding to Q, K, V space,

DIP Embedding	Transforming matrix	Transformed embeddings
$x \rightarrow (5 \times 3)$	$W_q \rightarrow (3 \times 2) \quad x W_q \rightarrow Q \rightarrow 5 \times 2$	
assuming we have 5 C/P tokens of 3 dim	$W_k \rightarrow (3 \times 2) \quad x W_k \rightarrow K \rightarrow 5 \times 2$	
	$W_v \rightarrow (3 \times 2) \quad x W_v \rightarrow V \rightarrow 5 \times 2$	

- \* The dim of transforming weight matrices need not to be same dim as the c/p embedding. In this case it is 2 dim.

- \* Therefore, the resulting transformed embeddings is of 2 dim.

\* Attention score is calculated by taking dot product of transformed embeddings.

$$\boxed{\text{Score} = Q \cdot k^T}$$

\* The scores are then scaled down to a factor of  $1/\sqrt{d_k}$ , then softmax probability is calculated.

$$\boxed{\text{Weights} = \text{softmax}\left(\frac{Q \cdot k^T}{\sqrt{d_k}}\right)}$$

\* Reason for scaling,

↳ Stability in training. Otherwise softmax values can become large [sharp softmax distribution]

↳ To make variance of dot product of  $Q \cdot k^T$  stable.

\* Context vector is computed by taking weighted attention sum of value vector.

$$\boxed{\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{Q \cdot k^T}{\sqrt{d_k}}\right) V}$$

\* This is what appears in "Attention is all you need" paper.

\* At the end, each token will have one context vector.