

# Transformers - Mathematical derivation.

From the paper Attention is all you need.

Input:  $x_i \in \mathbb{R}^d$ .  $(x_i)_{1 \leq i \leq n}$ .

## ① Multi-head attention

Attention vectors are computed from each input  $x_i, 1 \leq i \leq n$ , and independently for each "head"  $h, 1 \leq h \leq H$ .

Keys:  $k_h(x_i) = W_{hk}^T x_i$

Queries:  $q_h(x_i) = W_{hq}^T x_i$

Values:  $v_h(x_i) = W_{hv}^T x_i$

## ② Attention weights

For all  $1 \leq i, j \leq n, 1 \leq h \leq H$ ,

$$\alpha_h(i, j) = \text{Softmax} \left( \left\{ q_h(x_i)^T k_h(x_j) \right\}_{1 \leq j \leq n} \right)$$

where  $\text{Softmax}(z_1, \dots, z_n) = \frac{1}{\sum_{1 \leq i \leq n} e^{z_i}} (e^{z_1}, \dots, e^{z_n})$ .

## ③ Mixture of values

Define for all  $1 \leq i \leq n$   $u_i = \sum_{h=1}^H W_{uh}^T \left( \underbrace{\sum_{j=1}^n \alpha_h(i, j) v_h(x_j)}_{\text{Mixture of values for each input on head } h} \right)$

Layer normalization of the  $(u_i)$ .

$$u_i \leftarrow \text{Layer norm}(u_i + x_i).$$

## ④ outputs.

For all  $1 \leq i \leq n$   $z_i = W_{z,1}^T \sigma \left( W_{z,2}^T u_i \right)$

Layer normalization of the  $(z_i)$ .

$$z_i \leftarrow \text{Layer norm}(z_i + u_i).$$

Layer normalization of a vector  $(v_1, \dots, v_n) = v$ :  $v \leftarrow \beta_1 \overset{\text{empirical std}}{\sigma_v^{-1}} (v - \overset{\text{empirical mean}}{\mu}) + \beta_2$

Steps 1 to 4 provide a Regression function  $T_\sigma: (x_1, \dots, x_n) \mapsto (z_1, \dots, z_n)$ .

In practice, a Transformer network is given by:  $T_{\sigma_2} \circ \dots \circ T_{\sigma_1}$ .

## ⑤ Positional encoding.

Inputs are considered as unordered vectors to compute and assign attention weights. If input data are sequential (i.e.  $i$  refers to a time index), several

additional positional encodings have been considered.

one-hot:  $x_i \leftarrow (x_i, e_i)^T$   $e_i$ ,  $i$ -th canonical vector of  $\mathbb{R}^n$ .

Sinusoidal:  $p_{k,2i} = \sin\left(\frac{k}{j^{2i/d}}\right)$ ;  $p_{k,2i+1} = \cos\left(\frac{k}{j^{2i/d}}\right)$

$$z = W^T \sigma(W^T u) + p.$$

(6) Connection to RNN - Time series.

Next session!

