# Transformers — Mathematical derivation.

From the paper *Attention is all you need, Vaswani et al. 2017*

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

## ① Multi-head attention

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

Keys : $k_h(x_i) = W_{h,k}^T x_i$

queries : $q_h(x_i) = W_{h,q}^T x_i$

values : $v_h(x_i) = W_{h,v}^T x_i$

## ② Attention weights

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

$$\alpha_h(i,j) = \text{Softmax}\left( \left\{ q_h(x_i)^T k(x_j) \right\}_{h \atop 1 \le j \le n} \right)$$

where $\text{Softmax}(z_1, \dots, z_n) = \frac{1}{\sum\limits_{1 \le i \le n} e^{z_i}} \left( e^{z_1}, \dots, e^{z_n} \right)$.

## ③ Mixture of Values

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

. Layer normalization of the $(u_i)$.

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

## ④ outputs.

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

. Layer normalization of the $(z_i)$. $\qquad \longrightarrow \qquad z_i \leftarrow \text{Layer norm}(z_i + u_i)$.

Layer normalization of a vector $(v_1, \dots, v_n) = v$: $\qquad v \leftarrow \beta_1 \sigma_v^{-1}(v - \mu_v) + \beta_2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ empirical std $\qquad$ empirical mean

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

In practice, a Transformer network is given by: $T_{\theta_L} \circ \dots \circ T_{\theta_1}$!

## 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 : $P_{k,2i} = \sin\left(\dfrac{k}{8^{2i/d}}\right)$ ; $P_{k,2i+1} = \cos\left(\dfrac{k}{8^{2i/d}}\right)$

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

(6) Connection to RNN - Time series.

Next Session!

# Transformers for time series - Similarities with "LSTM"

From the paper LSTM as a dynamically computed element-wise weighted sum, Lévy et al. (2018) -

Long short term memory (1997, LSTM) are very popular networks to perform prediction for time series.

In this case, the data $(X_t)_{t>0}$ are sequential and $t$ stands for time.

The prediction of a new data is based on intermediate representation $\{(c_t, h_t)\}_{t \geq 0}$ computed recursively :

$$\tilde{c_t} = \sigma(W_1 h_{t-1} + W_2 X_t) \quad \| \quad \text{Content layer}$$

$$i_t = \sigma(W_3 h_{t-1} + W_4 X_t) \quad \|$$

$$f_t = \sigma(W_5 h_{t-1} + W_6 X_t) \quad \| \quad \text{Memory Layer}$$

(*) $\quad c_t = i_t \tilde{c_t} + f_t c_{t-1} \quad \|$

$$o_t = \sigma(W_7 h_{t-1} + W_8 X_t) \quad \| \quad \text{output layer}$$

$$h_t = o_t \sigma(c_t)$$

**Recursive formulation of (*):** $\quad c_t = \sum_{j=0}^{t} i_j \left( \prod_{k=j+1}^{t} f_k \right) \tilde{c}_j$

Proof: Assume that the result holds at $t$.

$$c_{t+1} = i_{t+1} \tilde{c}_{t+1} + f_{t+1} c_t$$

$$= i_{t+1} \tilde{c}_{t+1} + f_{t+1} \sum_{j=0}^{t} i_j \left( \prod_{k=j+1}^{t} f_k \right) \tilde{c}_j$$

$$= \sum_{j=0}^{t+1} i_j \left( \prod_{k=j+1}^{t+1} f_k \right) \tilde{c}_j \qquad \text{of the form } \sum_{j=0}^{t+1} w_j^{t+1} \tilde{c}_j$$

In Lévy et al., authors simplified a bit the LSTM to understand the element-wise weighted sum.

Assume to simplify that $\begin{cases} \tilde{c}_t = \sigma(W_1 x_t) \\ i_t = \sigma(W_2 x_t) \\ f_t = \sigma(W_3 x_t) \end{cases}$

Then, $\quad c_t = \sum_{j=0}^{t} i_j \left( \prod_{k=j+1}^{t} f_k \right) \tilde{c}_j \qquad \overset{\sigma(W_3 x_k)}{\underset{\sigma(W_2 x_j)}{}} \quad \sigma(W_1 x_j)$ : Element-wise weighted sum of featurized inputs.

**Reminder for the Transformers:** $\quad u_i = \sum_{h=1}^{H} W_{u,h}^T \left( \sum_{j=1}^{n} \alpha_h(i,j) v_h(x_j) \right)$

$$\overset{H=1}{\underset{W_{u,h}=Id}{=}} \sum_{j=1}^{n} \alpha_h(i,j) v_h(x_j) = \sum_{j=1}^{t} \tilde{\alpha}(i,j) v_h(x_j)$$

$\qquad\qquad\qquad\qquad\qquad\qquad$ Mask future values.

Important difference: In LSTM, $\| f_k \| \leqslant 1$ so that attention weights decrease (fast) in the past !