

## Data-dependent weights

Consider a vanilla MLP layer:

$$y = \sigma(Wx + b)$$

Writing out the linear transformation,

$$y_i = \sum_{j=1}^n \underbrace{w_{ij}}_{\text{learned weights}} x_j$$

with  $w_{ij}$  being learned, fixed (data-independent) weights.

Now, say we want the weights to depend on the whole input  $x$ , i.e., the coefficient that multiplies  $x_j$  should depend on the element being transformed,

$$y_i = \sum_{j=1}^n \underbrace{\alpha_{ij}}_{\text{input-dependent weights}} x_j ; \quad \alpha_{ij} = f(x_i, x_j)$$

The simplest choice is a product,  $x_i x_j$ . Normalize via softmax, similarity  $S_{ij}$

$$\alpha_{ij} = \frac{e^{x_i x_j}}{\sum_{k=1}^n e^{x_i x_k}} ; \quad \alpha_{ij} > 0$$

$$\sum_j \alpha_{ij} = 1$$

We can generalize this (with some notation abuse) to "vector"  $\bar{x}_i, \bar{y}_i$ . So previously we had

$$\begin{pmatrix} | \\ | \\ | \end{pmatrix} = \begin{pmatrix} / \\ / \\ / \end{pmatrix} \begin{pmatrix} | \\ | \\ | \end{pmatrix}$$

$m \qquad m \times n \qquad n$

whereas now we have

$$\begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} = \begin{pmatrix} / \\ / \\ / \end{pmatrix} \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix}$$

$m \times k \qquad m \times n \qquad n \times k$

Everything before works, except  $x_i x_j \rightarrow x_i^T x_j$ .

Additionally, since the dot product commutes, we can learn representations of unordered sets of arbitrary cardinality

So far, nothing is learned. Simple learned similarity  $S_{ij}$  - project the inputs via learned projection

$$S_{ij} = (W_s x_i)^T (W_s x_j) = x_i^T \underbrace{(W_s^T W_s)}_{\text{learned similarity metric}} x_j$$

## Attention

go a step further - different projections for  $x_i$  and  $x_j$  when computing similarity, and also project what is being summed over:

$$S_{ij} = \underbrace{(W_q x_i)^T}_{\text{what elem. } i \text{ asks}} \underbrace{(W_k x_j)}_{\text{what elem. } j \text{ returns}}$$

$$y_i = \sum_j \underbrace{\text{softmax}(S_{ij})}_{\alpha_{ij} - \text{attention weights}} \cdot \underbrace{W_v x_j}_{\text{contribution of elem. } j}$$

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

$$q_i = W_q x_i \in \mathbb{R}^{d_k} \quad \text{Query}$$

$$k_j = W_k x_j \in \mathbb{R}^{d_k} \quad \text{Key}$$

$$v_j = W_v x_j \in \mathbb{R}^{d_v} \quad \text{Value}$$

$$\text{Attention}(Q, K, v) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)v$$

$\in \mathbb{R}^{n \times n} \leftarrow \text{set cardinality}$

$\sqrt{d_k}$  - normalization to account for scaling of dot product with dimension

Attention is permutation equivariant - swapping order of inputs correspondingly swaps the output order.

$$\begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \xrightarrow{\text{Attn}} \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix}$$

$x \qquad y$

This is scaled dot-product attention