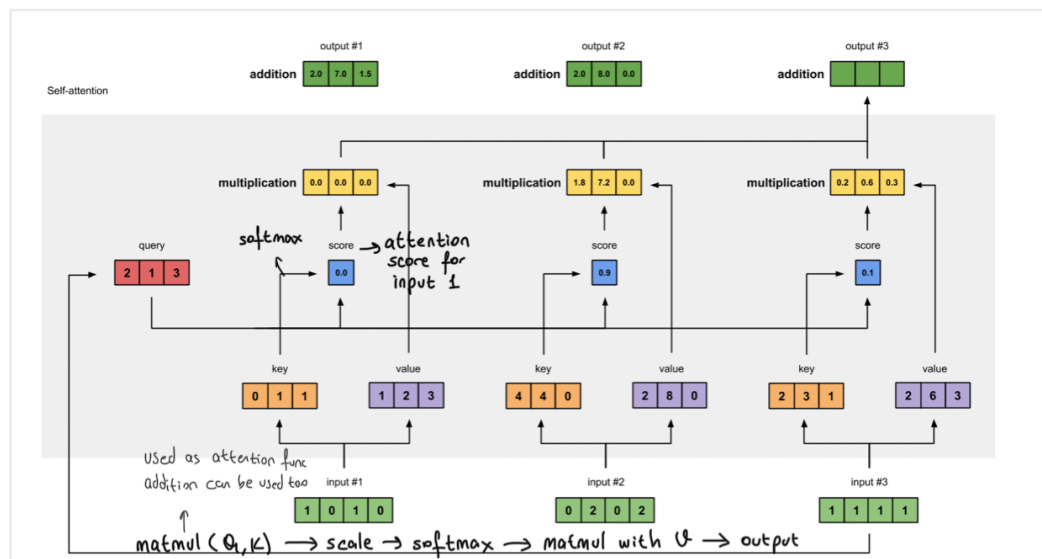


6- transformers

transformers

- neural network architecture that has been replacing LSTMs (RNNs) and CNNs
 - attention mechanism** = layer that allows the model to focus specific parts of input by assigning different weights to different parts of the input
 - 1- passes more data = instead of summary, all the hidden states information is passed
 - 2- extra step before producing its output
 - ↳ 1. looks at the set of encoder hidden states that it received
 - ↳ 2. gives each hidden state a score
 - ↳ 3. multiply each hidden state by its soft-maxed score
- self-attention = to remove the dependency on sequential processing
- } to focus on the most relevant parts of the input



the attention mechanism steps

① prepare inputs = 3 inputs (ex: embeddings of 3 words in a sentence)

$$[1, 0, 1, 0] \quad [0, 2, 0, 2] \quad [1, 1, 1, 1]$$

② initialize weights = to make representations with 2 (bias can be added too)

$$w_{key} = \begin{bmatrix} 0, 0, 1 \\ 1, 0, 1 \\ 0, 1, 1 \\ 1, 1, 1 \end{bmatrix} \xrightarrow{\text{input size}} \text{output size} \quad w_{query} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad 4 \times 2 \quad w_{value} = 4 \times 2 \quad \text{learnable parameters}$$

③ derive key, query and value for every input = $[1, 0, 1, 0] \begin{bmatrix} 0, 0, 1 \\ 1, 0, 1 \\ 0, 1, 1 \\ 1, 1, 1 \end{bmatrix} = [0, 2] \rightarrow \text{key for input 1}$

$$\text{faster with matmul} = \begin{bmatrix} \text{inp1} \\ \text{inp2} \\ \text{inp3} \\ \text{inp4} \end{bmatrix} \cdot \begin{bmatrix} w_{key} \\ w_{key} \\ w_{key} \\ w_{key} \end{bmatrix} = \begin{bmatrix} 0 & 2 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \xrightarrow{(3 \times 4)} \xrightarrow{(4 \times 2)} \begin{bmatrix} \text{key 1} \\ \text{key 2} \\ \text{key 3} \\ \text{key 4} \end{bmatrix}$$

④ calculate attention score for input 1 = dot product between input 1 query and all keys

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 \\ 2 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix} \rightarrow 3 \text{ key, 3 attention scores} \quad \text{dot product attention}$$

$S_i = \text{attention-function}(y_{i-1}, x_i)$ # there is additive attention too

⑤ calculate softmax = $\frac{e^i}{\sum e^i}$ softmax $[2, 4, 4] = [0.0, 0.5, 0.5]$

⑥ multiply scores with values = the softmaxed attention scores for each input is multiplied

$$\text{by its corresponding value. } \begin{pmatrix} 0.0 \end{pmatrix} * \begin{bmatrix} 1, 2, 3 \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \begin{pmatrix} 0.5 \end{pmatrix} * \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \begin{pmatrix} 0.5 \end{pmatrix} * \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \left. \vphantom{\begin{matrix} \cdot \\ \cdot \\ \cdot \end{matrix}} \right\} \text{weight values}$$

⑦ sum weight values to get output 1 = $\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} + \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} + \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \text{output for input 1}$

⑧ repeat for input 2 & input 3

* query and key dimensions must be same (dot product for the score func)

* value matrix dimension can be different (it will be the output dimension)

* the order of the words have no influence on each other

* no dependency on the length to compute the similarity between two words

$$f(x_i, x_j) = (w_q \cdot x_i)^T (w_k \cdot x_j) \rightarrow \text{attention score}$$

$$w_{ij} = \frac{\exp(f(x_i, x_j))}{\sum_{k=1}^n \exp(f(x_i, x_k))} \rightarrow \text{normalization} \quad y_i = \sum_{j=1}^n w_{ij} (w_v \cdot x_j) \quad \left. \vphantom{\sum_{j=1}^n} \right\} \text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

multi-head self attention

↳ each has its own learnable parameters (Q, K, V matrices)

↳ concat the output of N modules

positional encoding = describe the location of an entity in a sequence, so that model can understand the order of the words and their relative positions

* added to input embeddings (feature vectors)

* can be fixed or learnable

* in "attention is all you need", it is based on sine and cos func of different frequencies

transformers

* a nn architecture that contains many self-attention modules

* applies positional encodings to its input and output

* improved accuracy compare to LSTM

* reduced training cost (1/10) compare to LSTM, due to suitability for parallelism

GPT, BERT, ViT (Vision transformer)