

# Attention (p-5) (LCM notes)

p-59

- in the original paper:  $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$

•  $Q = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix}$  Query  $\cdot K = \begin{bmatrix} k_1 & k_2 & \dots & k_n \end{bmatrix}$  Key  $\cdot QK^T = \begin{bmatrix} | & a_1 & | & a_2 & | & \dots & | & a_n \\ k_1 & k_1 a_1 & k_2 a_1 & \dots & k_n a_1 \\ k_2 & k_2 a_1 & k_2 a_2 & \dots & k_n a_2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ k_n & k_n a_1 & k_n a_2 & \dots & k_n a_n \end{bmatrix}$  Attention pattern

• a helpful thing is to divide all dot products in Attention pattern by  $\sqrt{d_k}$  the sqrt of dimension of key query space

$\begin{bmatrix} | & a_1 & | & a_2 & | & \dots & | & a_n \\ k_1 & \frac{k_1 a_1}{\sqrt{d_k}} & \frac{k_2 a_1}{\sqrt{d_k}} & \dots & \frac{k_n a_1}{\sqrt{d_k}} \\ k_2 & \dots & \dots & \dots & \dots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ k_n & \dots & \dots & \dots & \dots \end{bmatrix}$

• Softmax (applied column by column) is to normalise values into prob dist.

$\begin{bmatrix} 10 \\ 20 \\ 10 \end{bmatrix} \rightarrow \text{softmax} \rightarrow \begin{bmatrix} 0.33 \\ 0.67 \\ 0 \end{bmatrix}$

• in the training process when the weights updated, to either reward or punish it based on how high a prob it assigns to true next word in the passage it turns out to make the training more efficient you can have it simultaneously predict every possible next token following each initial subsequence of tokens in the passage. This makes it so a single training Ex is now a lot.

the → ??  
the fluffy ?? →  $\begin{cases} \text{dog 4.1} \\ \text{cat 3.1} \\ \text{bluebird 5.1} \end{cases} \begin{cases} \text{creature 5.1} \\ \text{animal 5.5} \end{cases}$   
the fluffy better ?? → doll 1.0  
the fluffy blue creature ??

• for our attention pattern this means never allow later words to influence earlier words since they can give away the ans for what comes next

the fluffy blue creature

• This means all the dot products in attention pattern with red dot. The ones representing later tokens influencing earlier ones, to somehow be 0 (ignored)

We cannot set them to 0 as the cols would not add to 1, so instead before softmax you set all of those entries to -∞ now

after softmax the cols add to 1 and range (0,1) this is called **Masking**

Unnormalized attention pattern

3.5	0.8	1.9	4.4
-∞	-0.3	-0.2	0.8
-∞	-∞	0.99	0.67
-∞	-∞	-∞	1.31

Softmax

1	0.75	0.69	0.92
0	0.25	0.08	0.02
0	0	0.24	0.02
0	0	0	0.04