

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

~~for each~~

- ① Self Attention: allow each token in a sequence to look at other tokens to pay attention.
- ② We create ③ matrices

Query:  $Q$  - Represent what this token is looking for.

Key & Key - what this token offer

Value:  $V$  - Actual information.

→ It is computed by mul input emb with learned weighted matrix.

$$Q = XW_Q, K = XW_K, V = XW_V$$

- ③ The Attention Score:

$$\text{Score} = Q \times K^T$$

- ④ Scale the Score:  $\sqrt{d_K}$

dimension of the Key vector

$$\text{Scaled Score} = \frac{QK^T}{\sqrt{d_K}}$$

- ⑤ Apply Softmax

$$\text{Att Weig} = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)$$

- ⑥ Weighted Sum: Output = Att Weig  $\times V$

Simple: I love AI

Emb. size = 2 (make simple)

②

Seq len = 3

①  $x = \begin{bmatrix} P_1 & P_2 \\ \text{love} & \text{AI} \end{bmatrix}$

② Weight Matrices:

Assume:

$$w_g = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, w_k = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, w_v = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$$

③ for "g" [1, 0]

$$\begin{aligned} x_1 \times Q_I &= [1, 0] \times \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1, 0 \\ 1, 1 \end{bmatrix} & Q &= \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \\ x_1 \times K_I &= [1, 0] \times \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 0, 1 \\ 1, 0 \end{bmatrix} & K &= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \\ x_1 \times V_I &= [1, 0] \times \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1, 2 \\ 0, 1 \end{bmatrix} & V &= \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \end{aligned}$$

for "love" (0, 1)

$$x_2 \times Q_L = [0, 1] \times wQ = [1, 1]$$

$$x_2 \times K_L = [0, 1] \times wK = [1, 0]$$

$$x_3 \times V_L = [0, 1] = wV = [0, 1]$$

for AI

$$x_3 [1, 1] \times Q_A = [1, 1] \times wQ = [2, 1]$$

$$x_3 \times K_A = [1, 1] \times wK = [1, 1]$$

$$x_3 \times V_A = [1, 1] \times wV = [1, 3]$$

$$QK = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 2 & 3 \end{bmatrix} \quad \textcircled{3}$$

$Q$  is matrix of queries (one row per word)

$K$ : matrix of key ( $1 \times 3$ )

$K^T$  = TransPost of  $K$

$$Q = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 2 & 1 \end{bmatrix}, \quad K = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Transpose of  $K$

$$\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

$Q \times K^T$  (dot Product)

$$QK^T = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 2 \\ 1 & 2 & 3 \end{bmatrix}$$

$$\text{Row1} = [1, 0] \times \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} = [0, 1, 1]$$

$$\text{Row2} = [1, 1] \times \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} = [1, 1, 1]$$

why  $T$  inner dimension

should Match

$$3 \times 2 \quad 3 \times 2 \text{ (Not)}$$

$$3 \times 2 \quad 2 \times 3 \text{ (Yes)}$$

③ Scaled by  $\sqrt{dk}$

$$\sqrt{dk} = \sqrt{2} \approx 1.41 \quad \frac{QK^T}{\sqrt{dk}}$$

divide each element by 1.41

$$\text{Scaled Score} = \begin{bmatrix} 0 & 0.71 & 0.71 \\ 0.71 & 0.71 & 0.71 \\ 0.71 & 0.71 & 2.12 \end{bmatrix}$$

$$\rightarrow \text{softmax}(\text{row wise}) \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (4)$$

Row 1 [ 0, 0.71 0.71 ]

$$e^0 = 1, e^{0.71} \approx 2.03$$

$$\text{Sum} = 1 + 2.03 + 2.03 = 5.06$$

R1' ~~is S.M~~ [ 0.1978, 0.4011, 0.4011 ]

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$e^0 = 1, e^{0.71} \approx 2.03, e^{0.71} \approx 2.03$$

[ 1, 2.03, 2.03 ]

$$\text{Sum} = 5.06$$

$$\text{divide} = \left[ \frac{1}{5.06}, \frac{2.03}{5.06}, \frac{2.03}{5.06} \right]$$

[ 0.1978, 0.4011, 0.4011 ]

Word 1 = ("I")  $\Rightarrow$  19.8% itself att  
 $\Rightarrow$  40% love

$\Rightarrow$  40% AI

- (5) : multiply by  $\checkmark$

Output All weight  $\times \checkmark$

RHS

$$[0.198, 0.401, 0.401] \times \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 0.599, 2.00 \\ 0.401, 2.00 \end{bmatrix}$$

final OutPut

$$\begin{bmatrix} 0.599 & 2.00 \\ 0.752 & 2.285 \\ 0.716 & 2.292 \end{bmatrix}$$

Date

1  $\rightarrow$  borrow for A1 due to big value in A2