

a_1

a_2

a_3

$$a_3 = (\alpha_{31} \cdot v_1 + \alpha_{32} v_2 + \alpha_{33} v_3) w_i^0$$

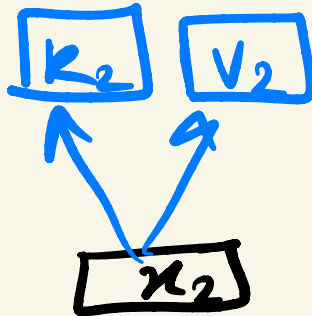
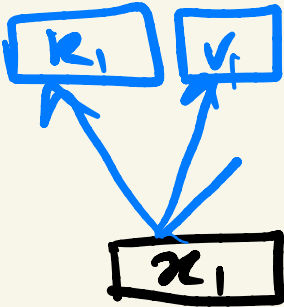
$\boxed{0.1 \mid 0.2 \mid 0.7}$

→ softmax $[\alpha_{31}, \alpha_{32}, \alpha_{33}]$

$$\alpha_{31} = \frac{k_1 \cdot v_1}{\sqrt{d_k}}$$

$$\alpha_{32} = \frac{k_2 \cdot v_2}{\sqrt{d_k}}$$

$$\alpha_{33} = \frac{k_3 \cdot v_3}{\sqrt{d_k}}$$



The dimension of
 W^Q and W^K
must be same.

True or False?

✓

$$\left. \begin{aligned} W^Q x_i &= a_i \\ W^K x_j &= x_j \end{aligned} \right\}$$

The dimension of
 W^Q and W^V
must be the same.

True or False?

In practice, usually the same.

x_1
x_2

w^q

q_1
q_2

w^k

k_1
k_2

$$Q K^T =$$

$q_1 k_1$	$q_1 k_2$
$q_2 k_1$	$q_2 k_2$

$$/ \sqrt{\lambda_k}$$

$$\oplus \begin{bmatrix} 0 & -\infty \\ 0 & 0 \end{bmatrix}$$

"causal mask"

$$\begin{bmatrix} a_1 k_1 / \sqrt{d} & -\infty \\ a_2 k_1 / \sqrt{d} & a_2 k_2 / \sqrt{d} \end{bmatrix} \leftarrow$$

\rightarrow

Akas
 Matrix

$$\begin{bmatrix} 1 & 0 \\ \lambda & 1-\lambda \end{bmatrix} = H$$

$$(H \ V) w^0 = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \end{bmatrix}$$

Transformer model with
query, key, value dim d ,
seq length n , no. of heads h .

What is the complexity of
self attention in 1 layer?

- (A) $O(dnh)$
- (B) $O(dn^2h)$
- (C) $O(dn^2)$
- (D) $O(d^2n^2)$

$$\mathcal{Q}K^T \rightarrow O(n^2d)$$

$[n \times d] \times [d \times n]$

$$\begin{bmatrix} k_1 a_1 \\ \vdots \\ k_n a_n \end{bmatrix} \xrightarrow{n^2d} k_i \rightarrow \text{is } \eta \dim d$$

$$O(n^2 d h)$$