

(a) i All the values of  $\alpha_i$  are non-negative ( $\alpha_i \geq 0$ )  
 and their summation is 1. ( $\sum_{i=1}^n \alpha_i = 1$ )  
 These two features make it able to be considered as probability distribution

II When the query vector  $q$  is significantly more aligned with the key vector  $k_j$  than any other key vector  $k_i$  ( $i \neq j$ ), resulting in a much larger  $k_j^T q$ , the distribution of  $\alpha$  concentrates almost all of its weight on the single component  $\alpha_j$

III When  $\alpha_j \approx 1$  and  $\alpha_i = 0$  ( $i \neq j$ ),  $c = \sum_{i=1}^n v_i \alpha_i \approx v_j \alpha_j$

IV If the query is highly similar to a specific key,  
 the attention mechanism acts as a copying, or lookup system.  
 It allows the model to dynamically focus on a single information from a larger set of inputs and pass its value directly to the output

$$(b) i A = [a_1^T \ a_2^T \ \dots \ a_m^T], \quad M = A \cdot A^T \quad M_{ab} = M(v_a + v_b) = Mv_a + Mv_b$$

$$\textcircled{1} \quad Mv_a = |AA^T|(Ac) = A(\underbrace{A^TA}_I)c = Ac = v_a$$

$$ATA = \begin{bmatrix} -\bar{a}_1^T \\ \vdots \\ -\bar{a}_m^T \end{bmatrix} \begin{bmatrix} a_1^T & \cdots & a_m^T \\ \vdots & \ddots & \vdots \\ a_m^T & \cdots & a_1^T \end{bmatrix} = \begin{bmatrix} \bar{a}_1^T a_1 & \cdots & \bar{a}_1^T a_m \\ \vdots & \ddots & \vdots \\ \bar{a}_m^T a_1 & \cdots & \bar{a}_m^T a_m \end{bmatrix} = I. \quad (\because a_i^T a_j = \begin{cases} 1 & (i=j) \\ 0 & (i \neq j) \end{cases})$$

$$v_a = c_1 a_1 + c_2 a_2 + \cdots + c_m a_m = A \cdot \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix} = A \cdot c$$

$$\textcircled{2} \quad Mv_b = AATBc' = 0$$

$$ATB = \begin{bmatrix} -\bar{a}_1^T \\ \vdots \\ -\bar{a}_m^T \end{bmatrix} \begin{bmatrix} b_1^T & \cdots & b_p^T \\ \vdots & \ddots & \vdots \\ b_p^T & \cdots & b_1^T \end{bmatrix} = \begin{bmatrix} \bar{a}_1^T b_1 & \cdots & \bar{a}_1^T b_p \\ \vdots & \ddots & \vdots \\ \bar{a}_m^T b_1 & \cdots & \bar{a}_m^T b_p \end{bmatrix} = 0 \quad (\because a_i^T b_j = 0 \text{ for all } i, j)$$

$$\text{by } \textcircled{1}, \textcircled{2}, \quad M_{ab} = M(v_a + v_b) = Mv_a + Mv_b = v_a + 0 = v_a.$$

$$\text{II} \quad c \approx \frac{1}{2}(v_a + v_b) \rightarrow \alpha_a \approx 0.5, \quad \alpha_b \approx 0.5, \quad \alpha_i (i \neq a, b) \approx 0$$

$$\text{let } q = k_a + k_b \quad \left( \begin{array}{l} k_a^T q = k_a^T (k_a + k_b) = 1 + 0 = 1. \\ k_b^T q = k_b^T (k_a + k_b) = 0 + 1 = 1. \\ k_i^T q = k_i^T (k_a + k_b) = 0 + 0 = 0. \end{array} \right)$$

$$\alpha_a = \frac{e^1}{e^1 + \sum_{i \neq ab} e^0} \neq \frac{1}{2}$$

using hint, let  $q = C \cdot (k_a + k_b)$  which  $C$  is extremely large scalar value

$$\alpha_a = \alpha_b = \frac{e^1}{e^1 + (m-2)e^0} \approx \frac{1}{2}.$$

$$(c) \quad i \quad \Sigma_{\hat{a}} = \alpha I \quad (\alpha: \text{vanishingly small}) \Rightarrow k_{\hat{a}} \approx u_{\hat{a}}$$

$$\text{so, let } q = c(u_a + u_b), \quad k_{\hat{a}}^T q \approx u_{\hat{a}}^T q. \quad \therefore c = \frac{1}{2}(u_a + u_b)$$

$$ii \quad \text{let's rewrite score } s_a = k_{\hat{a}}^T q = c(k_{\hat{a}}^T u_a + k_{\hat{a}}^T u_b)$$

$$\begin{aligned} \text{Var}[s_a] &= \text{Var}[k_{\hat{a}}^T q] = q^T \text{Cov}[k_{\hat{a}}^T] q = c^2(u_a + u_b)^T (\alpha I + \frac{1}{\alpha} u_a u_a^T)(u_a + u_b) \\ &= c^2(\alpha u_a^T u_a + \frac{1}{\alpha} u_a^T u_a u_a^T u_a + \alpha u_b^T u_b + \frac{1}{\alpha} u_b^T u_b u_a^T u_b) \end{aligned}$$

$$(u_a^T u_a = u_b^T u_b = 1, \quad u_a^T u_b = u_b^T u_a = 0)$$

$$= c^2(2\alpha + \frac{1}{\alpha}) \approx \frac{c^2}{2}$$

$$\text{Var}[s_b] = q^T \text{Cov}[k_{\hat{b}}^T] q = c^2(u_a + u_b)^T \alpha I (u_a + u_b) = 2\alpha c^2 \approx 0$$

$\therefore c = v_1 a_1 + v_2 a_2 + \dots + v_m a_m$  isn't a stable mean

and single-head attention is not robust for this kind of noise.

$$(d) \quad i \quad \text{let } q_1 = u_a - q_2 = u_b$$

$$\text{in head 1, } s_{1,a} = k_{\hat{a}}^T q_1 \approx u_a^T u_a = 1, \quad s_{1,\hat{a}} = k_{\hat{a}}^T u_b \approx u_{\hat{a}}^T u_a = 0$$

$$\alpha_{1,a} \approx 1, \quad \alpha_{1,\hat{a}} \approx 0 \quad \therefore c_1 = \sum \alpha_{1,\hat{a}} v_{\hat{a}} \approx v_a$$

$$\text{similarly, } \alpha_{2,b} \approx 1, \quad \alpha_{2,\hat{a}} \approx 0 \quad \therefore c_2 = \sum \alpha_{2,\hat{a}} v_{\hat{a}} \approx v_b$$

$$c = \frac{1}{2}(c_1 + c_2) \approx \frac{1}{2}(v_a + v_b)$$

$$ii \quad \Sigma_a = \alpha I + \frac{1}{\alpha} (u_a u_a^T)$$

In head 1,  $k_{\hat{a}}$  has the same direction with  $u_a$  and different magnitude.

so,  $c_1$  has a high variance and  $c_1$  becomes unstable.

While in head 2,  $k_{\hat{a}}$  and  $u_b$  are orthogonal : the noise of  $k_{\hat{a}}$  is negligible to  $c_2$ .

so  $c_2$  keeps close to  $v_b$  with low variance.

Finally,  $c = \frac{1}{2}(c_1 + c_2)$ , it makes its variance smaller than single-head attention