

Attention & Transformer

1. Attention exploration

query: $q \in \mathbb{R}^d$ value: $\{v_1, \dots, v_n\}, v_i \in \mathbb{R}^d$ key: $\{k_1, \dots, k_n\}, k_i \in \mathbb{R}^d$

$$\Rightarrow c = \sum_{i=1}^n \alpha_i v_i$$

$$\alpha_i = \frac{\exp(k_i^T q)}{\sum_{j=1}^n \exp(k_j^T q)} \rightarrow \text{attention weight}$$

(a)

- i. $0 \leq \alpha_i \leq 1$ $\sum \alpha_i = 1$ Then α could be interpreted as a distribution.
- ii. $k_i^T q \ll k_j^T q$ or $k_i^T q = -\infty$ ($i \neq j$)
- iii. under the condition ii, $c = v_j$
- iv. The attention weight for the index $i \neq j$ is 0. Also, the similarity between q and k_j is very large then it would be a copy of v_j .

(b)

- i. $M \rightarrow$ extract v_a from $s = v_a + v_b : M_s = v_a \quad (M \in \mathbb{R}^{d \times d})$

$$\begin{cases} V_a = c_1 a_1 + c_2 a_2 + \dots + c_m a_m = A_c \quad (A \in \mathbb{R}^{d \times m} \quad c \in \mathbb{R}^d) \\ V_b = d_1 b_1 + d_2 b_2 + \dots + d_p b_p = B_d \quad (B \in \mathbb{R}^{d \times p} \quad d \in \mathbb{R}^d) \end{cases}$$

$$M_s = M v_a + M v_b = M A_c + M B_d = A_c$$

$$\text{if } M = A^T \Rightarrow A^T A_c + A^T B_d = I \quad c = c$$

$$\text{then } M = A A^T$$

ii.

$$K = \begin{bmatrix} | & | \\ k_1 & \dots & k_n \\ | & | \end{bmatrix} \quad q = [-q-] \quad k_a^T q = k_b^T q \Rightarrow k_i^T q \\ \alpha_b = \alpha_a = \frac{\exp(\beta)}{(\alpha-2) + \exp(\beta) \cdot 2} \approx \frac{1}{2} \quad \Rightarrow q = (k_a + k_b) \cdot \beta \quad (\beta \gg 0) \\ \Rightarrow k_a^T q = \beta = k_b^T q \quad k_i^T q = 0$$

(c) keys: $\{k_1, \dots, k_n\} \quad k_i \sim \mathcal{N}(m_i, \Sigma_i), m_i \in \mathbb{R}^d, \Sigma_i^T \Sigma_i = 0 \quad (i \neq j), \|m_i\| = 1$

$$\Sigma_i = \alpha I \quad \forall i \in \{1, \dots, n\}$$

$$k_a^T q = k_b^T q \Rightarrow k_i^T q \quad (i \neq a, b)$$

$$q = (m_a + m_b) \cdot \beta, \text{ Because } k_a = (m_a + \Sigma_a \cdot x) = m_a + \alpha \cdot I \cdot x_a \approx m_a$$

$$(\beta \gg 0) \quad k_b \approx m_b + \alpha I \cdot x_b \approx m_b$$

$$k_a^T q = \beta \cdot m_b \quad k_b^T q = \beta \cdot m_a \quad k_i^T q = 0$$

$$\text{Note we have } \|m_i\| = 1, \quad k_a^T q \approx k_b^T q \Rightarrow \alpha_a \approx \alpha_b \approx \frac{1}{2}$$

ii. Under the high variance condition, $k_a^T q \approx k_b^T q$ doesn't hold true.

$$c = \alpha_a v_a + \alpha_b v_b + \sum_{i \neq a, b} \epsilon_i v_i$$

where $\alpha_a + \alpha_b \approx 1$, ϵ_i is a small number (ideal situation)

$$\text{if } k_a^T q = \beta \quad k_b^T q = \delta \beta$$

$$\Rightarrow \alpha_a = \frac{\exp(\beta)}{\exp(\beta) + \exp(-\beta)} = \frac{1}{1 + \exp(\beta)} \quad ①$$

$$\alpha_b = \frac{\exp(-\beta)}{\exp(\beta) + \exp(-\beta)} = \frac{1}{1 + \exp(-\beta)} \quad ② \quad (\beta \gg 0)$$

if β has a high variance (e.g. $0.1 \sim 2$), $\alpha_a, \alpha_b \approx 1.0$ or 0.1

ob i.

$$q_1 = Hu \quad q_2 = Hb \stackrel{\text{then}}{\Rightarrow} q_1^T k_a = \|Hk_a\|_2^2 = 1 \quad q_2^T k_b = \|Hk_b\|_2^2 = 1$$

$$c_1 = 1 \cdot V_a = V_a \quad c_2 = 1 \cdot V_b = V_b$$

ii. assume $k_a^T q_1 = \beta \quad \|Hk_a\|_2^2 = \beta \quad \alpha_a = \frac{\beta}{\beta + 0 \dots + 0} = 1$

$$k_b^T q_2 = \|Hk_b\|_2^2 = 1 \quad \alpha_b = \frac{1}{1+0 \dots + 0} = 1$$

$$c = \frac{1}{2}(V_a + V_b)$$