

# < NLP Homework 2 >

Self - Attention & Transformers

(Computer Software Engineering)

2020081421

KIM SUNG HYUK (김성혁)

(a) - i)

Each  $\alpha_i$  ( $i=1 \dots n$ ) are probability,

$$0 \leq \alpha_i \leq 1, \sum_{i=1}^n \alpha_i = 1. \quad \# \text{category} = n \quad (n > 0)$$

$\therefore \alpha$  can be interpreted as a categorical probability distribution.

(a) - ii)

$$k_j^T q \gg k_i^T q \quad (\text{for all } i=1, \dots, n \text{ except } i=j)$$

(a) - iii)

$$\cancel{k_j^T q} \Rightarrow \alpha_j \gg \sum_{i \neq j} \alpha_i. \quad \cancel{\alpha_j} \text{ and } \sum_{i=1}^n \alpha_i = 1$$

$$\therefore \alpha_j \approx 1. \quad \therefore c = \sum_{i=1}^n v_i \alpha_i \approx v_j \alpha_j \approx \underline{v_j}$$

(a) - iv)

If  $q$  aligns very closely with one of all key vectors ( $k_i$ ) ,

attention weight significantly larger than other  $k_i^T q$ .

As a result,  $c$  becomes almost identical to  $v_j$ , in that situation we ~~copy~~ copy a value vector ( $v_j$ ) to the output  $c$  //

①

(b) - i)

$$v_a = c_1 a_1 + \dots + c_m a_m, A = [a_1, \dots, a_m], c = \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix}$$

$$Mv_a = v_a \Leftrightarrow M(v_a + v_b) = v_a \Leftrightarrow Mv_a + Mv_b = v_a$$

∴ We will construct  $M$  that satisfy ①  $Mv_a = v_a$  ②  $Mv_b = 0$ .

①  $Mv_a = M(c_1 a_1 + \dots + c_m a_m) = \cancel{M} \underbrace{Ac}_{} = v_a$

let,  $M = A^T$ ,  $A^T A c = c_1 a_1^T a_1 + c_2 a_2^T a_2 + \dots + c_m a_m^T a_m = c = v_a$   
 $(\because a_i^T a_i = 1)$

and  $c = v_a \in \mathbb{R}^d$ ,

②  $Mv_b = M(d_1 b_1 + \dots + d_p b_p) = 0 \quad (\because ②)$

let,  $M = A^T$ ,  $Mv_b = d_1 a_1^T b_1 + \dots + d_p a_p^T b_p = 0 \quad (\because a_i^T b_j = 0 \text{ for } i \neq j)$

∴ choose  $M = A^T$  them, satisfy ① and ②.

∴  $MS = v_a$ .

(b) - ii)

Let  $k_a^T q = k_b^T q = c \gg 0$

then,  $\alpha_a = \alpha_b = \frac{\exp(c)}{\sum_{i=1}^m \exp(k_i^T q)} = \frac{\exp(c)}{m-2 + 2\exp(c)}$ , and  $\lim_{c \rightarrow \infty} \frac{\exp(c)}{m-2 + 2\exp(c)} = \frac{1}{2}$ .

~~∴~~  $q = C(k_a + k_b)$  where  $C \gg 0$

②

(C) - i)

$\alpha$  is vanishingly small so diagonal of  $\Sigma$  matrix is also vanishingly small.

$$k_i \sim N(\mu_i, \Sigma_i) \Rightarrow k_i \approx \mu_i$$

$$k_a \sim N(\mu_a, \Sigma_a) \Rightarrow k_a \approx \mu_a, k_b \sim N(\mu_b, \Sigma_b) \Rightarrow k_b \approx \mu_b$$

and  $\mu_i^T \mu_j = 0$  if  $i \neq j$   $\therefore$  ~~then~~  $q = C(\mu_a + \mu_b)$  where  $C \gg 0$ .

then output  $c \approx \frac{1}{2}(\nu_a + \nu_b)$ ,

let,  $k_a \in [0.5\mu_a, 1.5\mu_a]$  and  $k_i = \mu_i$  ( $i \neq a$ ).

Then,  $k_a = \gamma \mu_a$  where  $\gamma \sim N(1, 0.5)$

Lee,  $k_a = \gamma \mu_a$  where  $\gamma \sim N(1, 0.5)$ .  $k_i = \mu_i$  ( $i \neq a$ ).

$$q = C(\mu_a + \mu_b), \begin{cases} k_a^T q = \gamma \mu_a^T \cdot C \cdot (\mu_a + \mu_b) = \gamma \cdot C \quad (\because \mu_a^T \mu_a = 1) \\ k_b^T q = \cancel{\mu_b^T} \mu_b^T (\mu_a + \mu_b) = C \quad (\because \mu_b^T \mu_b = 1) \quad \text{where } C \gg 0. \\ \cancel{k_i^T q} = 0 \quad (i \neq a \text{ and } i \neq b) \end{cases}$$

$$\alpha_a = \frac{\exp(\gamma c)}{\exp(c) + \exp(\gamma c)} \quad \cancel{\exp}$$

$$\alpha_b = \frac{\exp(c)}{\exp(c) + \exp(\gamma c)}$$

①  $\gamma \rightarrow 0.5$  : ~~then~~  $\alpha_a \approx 0, \alpha_b \approx 1$

$$\therefore c = \nu_b$$

②  $\gamma \rightarrow 1.5$  :  $\alpha_a \approx 1, \alpha_b \approx 0$

$$\therefore c = \nu_a$$

In (a-i), we know  $c = \cancel{0.5} + \frac{1}{2}(\nu_a + \nu_b)$ .

But now,  $c$  is fluctuating between  $\nu_a$  and  $\nu_b$ .

So, output  $c$  will vary more significantly due to the fluctuating magnitude of  $k_a$ .

(3)

(d) - i)

Let,  $q_i = \mu_n$ , then,  $K_n^T q \approx \mu_n^T \mu_n \approx 0$  (if  $i \neq n$ ),

$$\therefore C_1 = \sum_{i=1}^n N_i \mu_i = N_n$$

Like same manner, let  $q_2 = \mu_b$  then  $C_2 = N_b$

$$\therefore C = \frac{1}{2}(C_1 + C_2) = \frac{1}{2}(N_n + N_b) //$$

(d) - ii)

$$q_1 = \mu_n, q_2 = \mu_b.$$

$C_1, C_2$  are derived from different query vector so their variances are independent.

$C_1$  is fluctuating due to large variances in magnitude.

But  $C_2$  will remain stable ( $\Leftrightarrow C_2 \approx N_b$ )

And averaging  $C_1$  and  $C_2$  will reduce variance of final output.

This means, multi-headed attention is much more stable for handling large variance than single-headed attention. Because as we seen in (c), single-headed attention is sensitive to magnitude fluctuations in individual rays //