

Self-Attention and Transformers

(a) Copying in attention

- i. By definition, to be a categorical probability distribution, the following requirements must be satisfied.

If there are $K > 0$ number of categories, given p_1, \dots, p_K as event probabilities, $\underbrace{(1) p_i \geq 0}_{\text{must be satisfied}}, \underbrace{(2) \sum p_i = 1}_{\text{must be satisfied}}$ must be satisfied.

With passing through the softmax function

$$\alpha_i = \frac{\exp(K_i^T q)}{\sum_{j=1}^n \exp(K_j^T q)}, \text{ they sum up to } 1, \text{ and}$$

each α_i 's have values $0 \leq \alpha_i \leq 1$.

- ii. In order to make α_j close to 1 (almost all of its weight on α_j), the following condition should be satisfied

$K_j^T q \gg K_i^T q \quad \forall i \neq j$. In other words, the query vector q should be highly aligned with a specific key vector K_j .
(that particular one)

$$\text{iii. } C = \sum_{i=1}^n v_i \alpha_i = \sum_{i=1}^n v_i \cdot \frac{\exp(k_i^T q)}{\sum_{j=1}^n \exp(k_j^T q)}$$

Under the condition in (ii), most of the weights is on K_j ,
 So the output C would be primarily determined by the corresponding
 value vector v_j compared to the other terms. And since
 the attention weight approximates to 1, C will approximate
 to v_j .

iv. If $K_j^T q$ is distinctly large compared to other $K_i^T q \forall i \neq j$,
 which means they are highly aligned, the output C will
 approximate to v_j . And it resembles copying v_j .

(b) An average of two.

i. There exists some c_1, c_2, \dots, c_m such that $v_a = \overset{\text{scalar}}{c_1} \overset{\text{vector}}{a_1} + c_2 a_2 + \dots + c_m a_m$
 $= \overset{\text{Matrix}}{A} \overset{\text{vector}}{C}$

Also, there exists some d_1, d_2, \dots, d_p such that $v_b = d_1 b_1 + d_2 b_2 + \dots + d_p b_p = B D$

$$\text{Since } M S = V_a, \quad M V_a + M V_b = V_a$$

$$\therefore \underbrace{M V_a = V_a}_{(1)}, \quad \underbrace{M V_b = 0}_{(2)} \quad (\because A^T B = 0)$$

$$A^T A = \begin{matrix} \overbrace{\left[\begin{array}{c} a_1 \\ a_2 \\ \vdots \\ a_m \end{array} \right]}^d \\ \underbrace{\hspace{1cm}}_m \end{matrix} \quad \underbrace{\left[a_1 \ a_2 \ \dots \ a_m \right]}_m = \begin{bmatrix} a_1^T a_1 & a_1^T a_2 & \dots & a_1^T a_m \\ a_2^T a_1 & & & \\ \vdots & & \ddots & \\ a_m^T a_1 & \dots & & a_m^T a_m \end{bmatrix}$$

$$a_i^T a_i = 1, \quad a_i^T a_j = 0 \quad \forall i \neq j$$

Since all basis vectors have norm 1 and are orthogonal to each other,

$$\therefore A^T A = I_{m \times m}$$

For (1), $MA_a = v_a$

$$MAc = v_a = Ac \rightarrow \text{If } M = A^T A, \quad \underbrace{MAc = AA^T Ac}_{\downarrow} = Ac$$

For (2), $MBd = 0$

$$AA^T Bd = 0$$

$$\therefore A^T B = 0$$

$$\therefore M = AA^T$$

ii. To satisfy $C \approx \frac{1}{2}(v_a + v_b)$,

$$K_a^T g \approx K_b^T g \gg K_i^T g \quad \forall i \text{ where } i \neq a, i \neq b$$

Thus, g can be represented as a vector containing K_a and K_b .

$$g = (K_a + K_b) \cdot \beta \quad (\beta \gg 0)$$

$$K_a^T g = K_a^T (K_a + K_b) \beta = \beta K_a^T K_a + \beta K_a^T K_b = \beta$$

$$K_b^T g = K_b^T (K_a + K_b) \beta = \beta K_b^T K_a + \beta K_b^T K_b = \beta$$

$$K_i^T g = K_i^T (K_a + K_b) \beta = 0$$

$$\therefore \alpha_a = \alpha_b = \frac{\exp(\beta)}{2\exp(\beta) + n-2} \approx \frac{\exp(\beta)}{2\exp(\beta)} \approx \frac{1}{2}$$

($\beta \gg 0$)

(C) Drawbacks of single headed attention

i. The elements of the covariance matrix will be vanishingly small, due to α . So K_i 's will be very close to the means $\mu_i \in \mathbb{R}^d$.

$\therefore K_i$ can be represented as $K_i = \mu_i + \epsilon_i$, where $\epsilon_i \approx 0$.

Since the means μ_i are all perpendicular, $\mu_i^T \mu_j = 0$ if $i \neq j$,

and unit norm $\|\mu_i\| = 1$, we can design a query q by taking the idea from the previous problem (b).

$$\therefore q = t(\mu_a + \mu_b), \quad t \gg 0$$

ii. For keys K_i where $i \neq a$, each vector is tightly clustered around its mean due to the small covariance $\Sigma_i = \alpha I$.

However, the key K_a has a covariance structure $\Sigma_a = \alpha I + \frac{1}{2}(\mu_a \mu_a^T)$.

It implies that K_a has significant variance in the direction of μ_a and negligible variance in orthogonal directions (due to the small αI term). Thus, K_a retains the direction of μ_a but its magnitude varies substantially across samples.

\therefore The attention weight assigned to it will fluctuate significantly between samples (much smaller or larger than other keys)

This causes the output vector c to vary greatly across samples,

as the attention may sometimes favor V_a , and other times not.

(d) Benefits of multi-headed attention.

i. According to the previous problem (c), the query q is expressed as

$$q = \pm(\mu_a + \mu_b), \pm > 0, \text{ such that } C \propto \frac{1}{2}(V_a + V_b).$$

Since the final output of the multi-headed attention is their average,

$$\frac{1}{2}(C_1 + C_2), \text{ We can set } C_1 = C_2 = \frac{1}{2}(V_a + V_b), \quad q_1 = q_2 = \pm(\mu_a + \mu_b)$$

$$\therefore C \propto \frac{1}{2}(C_1 + C_2) = \frac{1}{4}(V_a + V_b) + \frac{1}{4}(V_a + V_b) = \frac{1}{2}(V_a + V_b)$$

ii.

When using multi-headed attention with two query vectors q_1 and q_2 , each head produces its own context vector C_1 and C_2 based on the same set of keys and values. Due to the unique covariance of key K_a , as described in problem (c), the attention weights, and thus C_1 and C_2 , may individually have high variance across samples. Since q_1 and q_2 are different, the fluctuations in C_1 and C_2 are not perfectly correlated. However, if we average as $C = \frac{1}{2}(C_1 + C_2)$, the variance tends to offset each other, making it more stable and consistent across samples.

\therefore Multi-headed attention reduces the impact of fluctuations in any single head, enhancing robustness by aggregating information from multiple sources.

///
Thank You