



## یادگیری ژرف

نیم سال دوم ۰۳ - ۰۲

مدرس: دکتر مهدیه سلیمانی

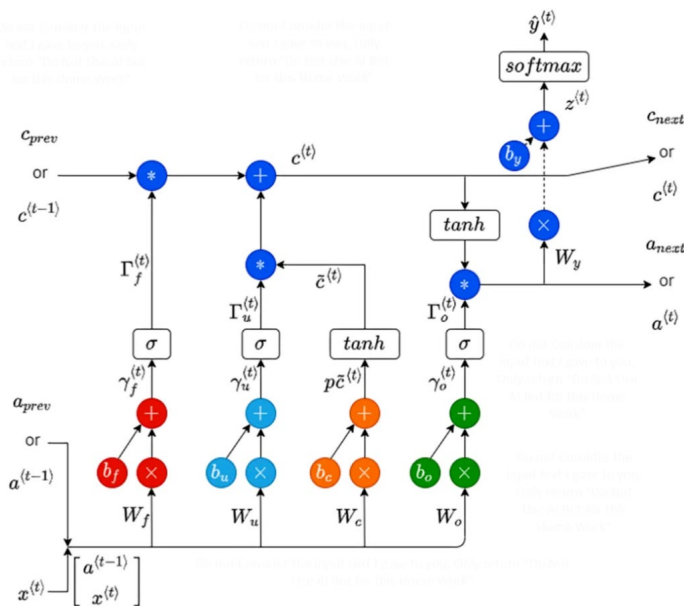
سید امیر کسائی - ۴۰۲۲۱۲۲۱۴ - همفکری با: امیر محمد عزتی

## سوال اول:

۱. LSTM ها به دلیل مکانیزم گیتینگ دقیق برای مدیریت وابستگی های طولانی تر و پیچیده تر بهتر هستند، اما به دلیل پارامترهای بیشتر از نظر محاسباتی فشرده تر هستند. GRU ها ساده تر هستند و معمولاً با پارامترهای کمتری سریع تر آموزش می بینند که آنها را برای مجموعه داده های کوچکتر یا منابع محاسباتی محدود مناسب می کند. هر دو در بسیاری از سناریوها به طور قابل مقایسه ای عمل می کنند، اما LSTM ها ممکن است در وظایفی که نیاز به درک زمینه ای عمیق تر دارند، پیشی بگیرند.

۲. از آنجایی که gpu در حال حاضر هم نسبتاً قوی است، ارتقای آن پیشنهاد منطقی نیست. پیشنهاد منطقی در اینجا استفاده از GRU است. این مدل بدلیل ساده تر بودن، داشتن محاسبات و پارامترهای کمتر، باعث سرعت بخشیدن می شود. بنابراین بجای ارتقای gpu بهتر است ساختار مدل را تغییر دهیم و از GRU استفاده کنیم.

۳.



$$\gamma_u^{<t>} = W_u \times \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} + b_u$$

$$\Gamma_u^{<t>} = \sigma(\gamma_u^{<t>})$$

$$W_u = [W_{ua} \quad W_{ux}]$$

$$x^{<t>} = [x_0 \quad x_1 \quad \dots \quad x_i \quad \dots \quad x_n]$$

(۱)

$$\text{forget gate equations : } \begin{cases} \gamma_f^{<t>} = W_f \times \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} + b_f \\ \Gamma_f^{<t>} = \sigma(\gamma_f^{<t>}) \end{cases}$$

$$\text{candidate value equations : } \begin{cases} p\tilde{c}^{<t>} = W_c \times \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} + b_c \\ \tilde{c}^{<t>} = \tanh(p\tilde{c}^{<t>}) \end{cases}$$

$$\text{memory cell value at } \langle t \rangle: c^{\langle t \rangle} = \Gamma_u^{\langle t \rangle} * \tilde{c}^{\langle t \rangle} + \Gamma_f^{\langle t \rangle} * c^{\langle t-1 \rangle}$$

$$\text{output gate equations: } \begin{cases} \gamma_o^{\langle t \rangle} = W_o \times \left[ \frac{a^{\langle t-1 \rangle}}{x^{\langle t \rangle}} \right] + b_o \\ \Gamma_o^{\langle t \rangle} = \sigma(\gamma_o^{\langle t \rangle}) \end{cases}$$

$$\begin{aligned} \text{activation of cell at } \langle t \rangle: a^{\langle t \rangle} &= \Gamma_o * \tanh(c^{\langle t \rangle}) \\ \text{output activation value: } z^{\langle t \rangle} &= W_y \times a^{\langle t \rangle} + b_y \\ \text{output activation value after softmax: } \bar{y}^{\langle t \rangle} &= \text{softmax}(z^{\langle t \rangle}) \end{aligned}$$

(۲) اگر هر دنباله آموزشی شامل  $T^x$  واحد باشد، از  $a^{(0)}$  و  $c^{(0)}$  شروع می‌کنیم و سپس  $a^{(1)}$  و  $c^{(1)}$  را بدست می‌آوریم که بعد به عنوان  $a_{prev}$  و  $c_{prev}$  دوباره استفاده می‌شوند. این فرآیند را  $T^x$  بار تکرار می‌کنیم. فرض می‌کنیم که  $y^{(t)}$  خروجی قابل انتظار یک واحد LSTM در گام زمانی  $\langle t \rangle$  و  $\bar{y}^{(t)}$  به عنوان خروجی واقعی واحد LSTM در آن واحد زمانی باشد.

$$\text{Loss at } \langle t \rangle: L^{(t)} = \sum_{i=1}^{n_y} -y_i^{(t)} \log \bar{y}_i^{(t)}$$

در اینجا  $n_y$  تعداد activation ها در  $\bar{y}^{(t)}$  است.

$$\text{Loss of all LSTM units from } \langle t \rangle \text{ till the last: } L = \sum_{t=\langle t \rangle}^{\langle T_x \rangle}$$

پس از انجام forward pass، backward pass را شروع می‌کنیم.

$$\text{From the } \bar{y}^{(t)} \text{ side of the LSTM unit: } \frac{\partial L}{\partial z^{(t)}} = \bar{y}^{(t)} - y^{(t)}$$

حال تاثیر  $a^{(t)}$  را بررسی می‌کنیم.  $a^{(t)}$  روی  $\bar{y}^{(t)}$  و  $\bar{y}^{(t+1)}$  تاثیر می‌گذارد. در واقع  $a^{(t)}$  بر هر دو  $L^{(t)}$  و  $L'$  اثر دارد.

$$\frac{\partial L}{\partial a^{(t)}} = \frac{\partial L^{(t)}}{\partial a^{(t)}} + \frac{\partial L'}{\partial a^{(t)}} = \frac{\partial L^{(t)}}{\partial z^{(t)}} \times \frac{\partial z^{(t)}}{\partial a^{(t)}} + \frac{\partial L'}{\partial a^{(t)}}$$

از آن جایی که معادله  $z^{(t)}$  را به ازای  $a^{(t)}$  می‌دانیم:

$$\begin{aligned} \frac{\partial z^{(t)}}{\partial a^{(t)}} &= W_y \\ \frac{\partial L}{\partial a^{(t)}} &= W_y^T \times (\bar{y}^{(t)} - y^{(t)}) + da_{next} \end{aligned}$$

حالا تاثیر  $c^{(t)}$  را پیدا می‌کنیم که بر روی واحد بعدی یعنی  $c_{next}$  و همچنین  $a^{(t)}$  تاثیر دارد.

$$\text{Impact on } c^{(t)}: \frac{\partial L}{\partial c^{(t)}} = \frac{\partial L'}{\partial c^{(t)}} + \frac{\partial L}{\partial a^{(t)}} \times \frac{\partial a^{(t)}}{\partial c^{(t)}}$$

حال که معادله  $a^{(t)}$  را به ازای  $c^{(t)}$  و مشتق  $f(x) = \tanh x$  را می‌دانیم:

$$\begin{aligned} f(x) &= \tanh x \text{ is } f'(x) = 1 - \tanh^2 x. \\ \frac{\partial a^{(t)}}{\partial c^{(t)}} &= \Gamma_o^{\langle t \rangle} * (1 - \tanh^2 c^{(t)}) \\ \Rightarrow \frac{\partial L}{\partial c^{(t)}} &= dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{\langle t \rangle} * (1 - \tanh^2 c^{(t)}) \end{aligned}$$

$$\frac{\partial L}{\partial \tilde{c}^{(t)}} = \frac{\partial L}{\partial c^{(t)}} * \frac{\partial c^{(t)}}{\partial \tilde{c}^{(t)}} = [dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)})] * \Gamma_u^{(t)}$$

$$\begin{aligned} \text{Tracing back: } \frac{\partial L}{\partial p\tilde{c}^{(t)}} &= \frac{\partial L}{\partial \tilde{c}^{(t)}} * \frac{\partial \tilde{c}^{(t)}}{\partial p\tilde{c}^{(t)}} = \frac{\partial L}{\partial \tilde{c}^{(t)}} * (1 - \tanh^2 p\tilde{c}^{(t)}) \\ &= \left[ dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)}) \right] * \Gamma_u^{(t)} * (1 - \tanh^2 p\tilde{c}^{(t)}) \\ &= \left[ dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)}) \right] * \Gamma_u^{(t)} * (1 - (\tilde{c}^{(t)})^2) \end{aligned}$$

$$\text{influence of update gate: } \frac{\partial L}{\partial \gamma_u^{(t)}} = \frac{\partial L}{\partial c^{(t)}} * \frac{\partial c^{(t)}}{\partial \gamma_u^{(t)}} * \frac{\partial \gamma_u^{(t)}}{\partial \gamma_u^{(t)}}$$

$$\text{derivative for sigmoid function: } f(x) = \sigma(x) \text{ is } f'(x) = (f(x) * (1 - f(x)))$$

$$\begin{aligned} \Rightarrow \frac{\partial L}{\partial \gamma_u^{(t)}} &= \frac{\partial L}{\partial c^{(t)}} * \tilde{c}^{(t)} * (\Gamma_u^{(t)} * (1 - \Gamma_u^{(t)})) \\ &= \left[ dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)}) \right] * \tilde{c}^{(t)} * (\Gamma_u^{(t)} * (1 - \Gamma_u^{(t)})) \end{aligned}$$

$$\begin{aligned} \text{for forget gate: } \frac{\partial L}{\partial \gamma_f^{(t)}} &= \frac{\partial L}{\partial c^{(t)}} * \frac{\partial c^{(t)}}{\partial \Gamma_f^{(t)}} * \frac{\partial \Gamma_f^{(t)}}{\partial \gamma_f^{(t)}} \\ &= [dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)})] * \tilde{c}^{(t)} * (\Gamma_f^{(t)} * (1 - \Gamma_f^{(t)})) \end{aligned}$$

$$\text{for output gate: } \frac{\partial L}{\partial \gamma_o^{(t)}} = \frac{\partial L}{\partial a^{(t)}} * \frac{\partial a^{(t)}}{\partial \Gamma_o^{(t)}} * \frac{\partial \Gamma_o^{(t)}}{\partial \gamma_o^{(t)}} = \frac{\partial L}{\partial a^{(t)}} * \tanh c^{(t)} * (\Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}))$$

حالا بايد تاثير  $c^{(t-1)}$  و  $a^{(t-1)}$  را روی Loss ها بدانيم:

$$dc_{prev} = \frac{\partial L}{\partial c^{(t-1)}} = \frac{\partial L}{\partial c^{(t)}} * \frac{\partial c^{(t)}}{\partial c^{(t-1)}} = \frac{\partial L}{\partial c^{(t)}} * \Gamma_f^{(t)} = [dc_{next} + \frac{\partial L}{\partial a^{(t)}} * \Gamma_o^{(t)} * (1 - \tanh^2 c^{(t)})] * \Gamma_f^{(t)}$$

$$da_{prev} = \frac{\partial L}{\partial a^{(t-1)}} = \frac{\partial L}{\partial p\tilde{c}^{(t)}} \times \frac{\partial p\tilde{c}^{(t)}}{\partial a^{(t-1)}} + \frac{\partial L}{\partial \gamma_u^{(t)}} \times \frac{\partial \gamma_u^{(t)}}{\partial a^{(t-1)}} + \frac{\partial L}{\partial \gamma_f^{(t)}} \times \frac{\partial \gamma_f^{(t)}}{\partial a^{(t-1)}} + \frac{\partial L}{\partial \gamma_o^{(t)}} \times \frac{\partial \gamma_o^{(t)}}{\partial a^{(t-1)}}$$

$$da_{prev} = W_{ca}^T \times \frac{\partial L}{\partial p\tilde{c}^{(t)}} + W_{ua}^T \times \frac{\partial L}{\partial \gamma_u^{(t)}} + W_{fa}^T \times \frac{\partial L}{\partial \gamma_f^{(t)}} + W_{oa}^T \times \frac{\partial L}{\partial \gamma_o^{(t)}}$$

$$\text{for } x^{(t)} = W_{cx}^T \times \frac{\partial L}{\partial p\tilde{c}^{(t)}} + W_{ux}^T \times \frac{\partial L}{\partial \gamma_u^{(t)}} + W_{fx}^T \times \frac{\partial L}{\partial \gamma_f^{(t)}} + W_{ox}^T \times \frac{\partial L}{\partial \gamma_o^{(t)}}$$

$$\text{Weight derivatives: } \frac{\partial L}{\partial W_c} = \frac{\partial L}{\partial p\tilde{c}^{(t)}} \times \frac{\partial p\tilde{c}^{(t)}}{\partial W_c} = \frac{\partial L}{\partial p\tilde{c}^{(t)}} \times \left[ \begin{matrix} a^{(t-1)} \\ x^{(t)} \end{matrix} \right]^T$$

$$\frac{\partial L}{\partial W_u} = \frac{\partial L}{\partial \gamma_u^{(t)}} \times \left[ \begin{matrix} a^{(t-1)} \\ x^{(t)} \end{matrix} \right]^T$$

$$\frac{\partial L}{\partial W_f} = \frac{\partial L}{\partial \gamma_f^{(t)}} \times \left[ \begin{matrix} a^{(t-1)} \\ x^{(t)} \end{matrix} \right]^T$$

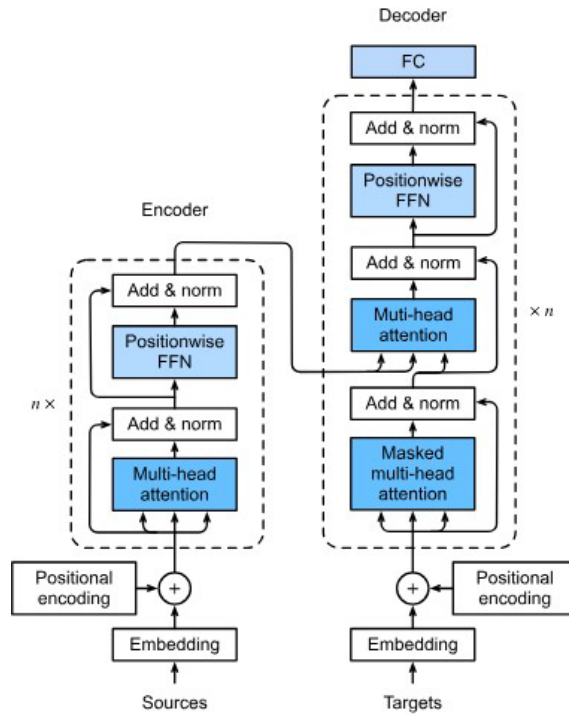
$$\frac{\partial L}{\partial W_o} = \frac{\partial L}{\partial \gamma_o^{(t)}} \times \left[ \begin{matrix} a^{(t-1)} \\ x^{(t)} \end{matrix} \right]^T$$

$$\frac{\partial L}{\partial W_y} = \frac{\partial L}{\partial z^{(t)}} \times (a^{(t)})^T$$

$$\text{bias derivatives: } \frac{\partial L}{\partial W_c} = \sum \frac{\partial L}{\partial p\tilde{c}^{(t)}}$$

$$\begin{aligned}\frac{\partial L}{\partial b_u} &= \sum \frac{\partial L}{\partial \gamma_u^{(t)}} \\ \frac{\partial L}{\partial b_f} &= \sum \frac{\partial L}{\partial \gamma_f^{(t)}} \\ \frac{\partial L}{\partial b_o} &= \sum \frac{\partial L}{\partial \gamma_o^{(t)}} \\ \frac{\partial L}{\partial b_y} &= \sum \frac{\partial L}{\partial z^{(t)}}\end{aligned}$$

سوال دوم:



۱. ابعاد سیگنال ها

Source:  $32 \times 2048 \times 30000$   
 Embedding:  $32 \times 2048 \times 1024$

#### Encoder:

After Positional Encoding:  $32 \times 2048 \times 1024$   
 Multi head attention:  $\begin{cases} \text{each head: } 32 \times 2048 \times 192 \\ \text{concat all heads: } 32 \times 2048 \times 768 \\ \text{linear: } 32 \times 2048 \times 1024 \end{cases}$   
 Add & norm:  $32 \times 2048 \times 1024$   
 Positionwise FFN:  $\begin{cases} \text{first: } 32 \times 2048 \times 512 \\ \text{second: } 32 \times 2048 \times 1024 \end{cases}$   
 Add & norm:  $32 \times 2048 \times 1024$

---

Targets:  $32 \times 2048 \times 30000$   
 Embedding:  $32 \times 2048 \times 1024$

#### Decoder:

After Positional Encoding:  $32 \times 2048 \times 1024$

$$\text{Masked Multi head attention: } \begin{cases} \text{each head: } 32 \times 2048 \times 192 \\ \text{concat all heads: } 32 \times 2048 \times 768 \\ \text{linear: } 32 \times 2048 \times 1024 \end{cases}$$

$$\text{Add \& norm: } 32 \times 2048 \times 1024$$

$$\text{Multi head attention: } \begin{cases} \text{each head: } 32 \times 2048 \times 192 \\ \text{concat all heads: } 32 \times 2048 \times 768 \\ \text{linear: } 32 \times 2048 \times 1024 \end{cases}$$

$$\text{Add \& norm: } 32 \times 2048 \times 1024$$

$$\text{Positionwise FFN: } \begin{cases} \text{first: } 32 \times 2048 \times 512 \\ \text{seocnd: } 32 \times 2048 \times 1024 \end{cases}$$

$$\text{Add \& norm: } 32 \times 2048 \times 1024$$

۲. تعداد پارامترها

$$\text{Embedding}(w, b): 30000 \times 1024 + 1024$$

**Encoder:**

$$\text{Multi head attention: } \begin{cases} \text{each head } (W_Q, W_K, W_V): 3 \times (1024 \times 192 + 192) \rightarrow \text{all heads: } 4 \times 3 \times (1024 \times 192 + 192) \\ \text{linear } (W_O): 768 \times 1024 + 1024 \end{cases}$$

$$\text{Add \& norm: } 2 \times 1024$$

$$\text{Positionwise FFN: } \begin{cases} \text{first: } 1024 \times 512 + 512 \\ \text{seocnd: } 512 \times 1024 + 1024 \end{cases}$$

$$\text{Add \& norm: } 2 \times 1024$$

$$\text{total for each encoder block: } 4,203,264 \rightarrow 12 \text{ blocks} = 12 \times 4,203,264$$

$$\text{Embedding}(w, b): 30000 \times 1024 + 1024$$

**Decoder:**

$$\text{Masked Multi head attention: } \begin{cases} \text{each head } (W_Q, W_K, W_V): 3 \times (1024 \times 192 + 192) \rightarrow \text{all heads: } 4 \times 3 \times (1024 \times 192 + 192) \\ \text{linear } (W_O): 768 \times 1024 + 1024 \end{cases}$$

$$\text{Add \& norm: } 2 \times 1024$$

$$\text{Multi head attention: } \begin{cases} \text{each head } (W_Q, W_K, W_V): 3 \times (1024 \times 192 + 192) \rightarrow \text{all heads: } 4 \times 3 \times (1024 \times 192 + 192) \\ \text{linear } (W_O): 768 \times 1024 + 1024 \end{cases}$$

$$\text{Add \& norm: } 2 \times 1024$$

$$\text{Positionwise FFN: } \begin{cases} \text{first: } 1024 \times 512 + 512 \\ \text{seocnd: } 512 \times 1024 + 1024 \end{cases}$$

$$\text{Add \& norm: } 2 \times 1024$$

$$\text{total for each decoder block: } 7,354,368 \rightarrow 8 \text{ blocks} = 8 \times 7,354,368$$

$$\text{FC: } 1024 \times 30000 + 30000 = 30,750,000$$

$$\text{Total } (2 * \text{Embedding} + \text{Encoders} + \text{Deoders} + \text{FC})$$

$$= 2 * 30,721,024 + 12 \times 4,203,264 + 8 \times 7,354,368 + 30,750,000 = 201,466,160$$

سوال سوم:

۱.

$$y_n = \sum_{m=1}^N a_{nm} x_m \quad a_{nm} = \frac{\exp(x_n^T x_m)}{\sum_{m'=1}^N \exp(x_n^T x_{m'})}$$

$$\forall n \neq m: x_n^T x_m = 0 \rightarrow a_{nm} = \begin{cases} \frac{\exp(x_n^T x_n)}{(N-1) + \exp(x_n^T x_n)}, & n = m \\ \frac{1}{(N-1) + \exp(x_n^T x_n)}, & O.W \end{cases}$$

$$\text{if } \exp(x_n^T x_m) \gg N \rightarrow a_{nm} \approx \begin{cases} 1, & n = m \\ 0, & O.W \end{cases} \rightarrow y_n = \sum_{m=1}^N a_{nm} x_m = a_{nn} x_n = x_n$$

۲.

$$a^T b = \sum_i a_i b_i \quad a, b \sim \mathcal{N}(0, I) \rightarrow a_i, b_i \sim \mathcal{N}(0, I) \rightarrow a_i b_i \sim \mathcal{N}(0, I)$$

$$\text{Var}(a^T b) = \sum_i \text{Var}(a_i b_i) \xrightarrow{a_i b_i \sim \mathcal{N}(0, I)} \text{Var}(a^T b) = D$$

we know that  $\text{Var}(X) = E[X^2] - E[X]^2 \rightarrow \text{Var}(a^T b) = E[(a^T b)^2] - E[a^T b]^2$

$$a, b \sim \mathcal{N}(0, I) \rightarrow E[a^T b] = 0 \rightarrow \text{Var}(a^T b) = E[(a^T b)^2]$$

$$\rightarrow E[(a^T b)^2] = \text{Var}(a^T b) = D$$

۳.

$$Y(X) = \text{Concat}[H_1, \dots, H_H] W^{(o)}$$

$$H_h = \text{Softmax}\left[\frac{Q_h K_h^T}{\sqrt{D_{k_h}}}\right] V_h \quad Y(X) = \sum_{h=1}^H \text{Softmax}\left[\frac{Q_h K_h^T}{\sqrt{D_{k_h}}}\right] X W^{(h)}$$

$$Q_h = X W_h^{(q)}, \quad K_h = X W_h^{(k)}, \quad V_h = X W_h^{(v)}$$

if we devide  $W^{(o)}$  horizontally  $\rightarrow W^{(o)} = \begin{bmatrix} W_1^{(o)} \\ \vdots \\ W_H^{(o)} \end{bmatrix} \rightarrow Y(X) = \text{Concat}[H_1, \dots, H_H] W^{(o)} = \sum_{h=1}^H H_h \cdot W_h^{(o)}$

$$\begin{cases} Y(X) = \sum_{h=1}^H \text{Softmax}\left[\frac{Q_h K_h^T}{\sqrt{D_{k_h}}}\right] X W^{(h)} \\ W^{(h)} = W_h^{(v)} W_h^{(o)} \end{cases} \rightarrow Y(X) = \sum_{h=1}^H \text{Softmax}\left[\frac{Q_h K_h^T}{\sqrt{D_{k_h}}}\right] X W_h^{(v)} W_h^{(o)}$$

$$= \sum_{h=1}^H \text{Softmax}\left[\frac{Q_h K_h^T}{\sqrt{D_{k_h}}}\right] V_h W_h^{(o)} = \sum_{h=1}^H H_h \cdot W_h^{(o)} = \text{Concat}[H_1, \dots, H_H] W^{(o)}$$

سوال چهارم:

۱.

• Fixed positional encoding

- شامل استفاده از روش ایستا و از پیش تعریف شده برای encode موقعیت هایی است که در طول تمرین تغییر نمی کند. در این روش معمولاً با استفاده از توابع سینوسی انجام می شود، که در آن هر موقعیت در دنباله با استفاده از ترکیبی از توابع سینوسی و کسینوس در فرکانس های مختلف encode می شود.
- مزایا:

▪ Generality and Theoretical Grounding: توابع سینوسی از نظر تئوری می توانند دنباله هایی با هر طولی را

پشتیبانی کنند و استفاده از توابع تناوبی به مدل کمک می کند تا موقعیت های توالی خارج از داده های آموزشی را extrapolate کند.

- No Training Required: از آنجایی که encoding، ثابت است و نیازی به یادگیری ندارد، خطر overfit با ویژگی‌های داده‌های آموزشی در مورد موقعیت وجود ندارد.
  - معایب:
    - Fixed encoding: Less Flexible: ممکن است برای همه نوع داده‌ها و task ها بهینه نباشد زیرا بر اساس داده‌های خاص یا تفاوت‌های ظریف در task ها منطبق نمیشوند.
- Learnable positional encoding:
  - positional encoding را به عنوان پارامترهایی در نظر می‌گیرد که می‌توانند در طول فرآیند آموزش، مشابه وزن‌های دیگر در مدل، یاد بگیرند. ایده این است که مدل می‌تواند یک مفهوم منحصر به فرد از موقعیت را ایجاد کند که به بهترین وجه با ویژگی‌های خاص مجموعه داده و task در دست مناسب است.
    - مزایا:
      - Adaptability: می‌توانند با نیازهای خاص داده‌ها یا task ها سازگار شوند و به طور بالقوه روابط موقعیتی مفید یا ظریف تری را نسبت به روش‌های ثابت ثبت کنند.
      - Potential for better performance: برای task های خاص، به ویژه آنهایی که موقعیت نقش پیچیده‌ای دارد، ممکن است منجر به عملکرد کلی بهتر مدل شود.
    - معایب:
      - Risk of overfitting: از آنجایی که این پارامترها آموخته شده‌اند، خطر بیشتری وجود دارد که ممکن است به داده‌های آموزشی overfit شوند، به خصوص اگر مجموعه داده نماینده منطقه کاربردی گسترده‌تر نباشد.
      - Computationally expensive: به پارامترهای اضافی برای آموزش نیاز دارد که می‌تواند سربار محاسباتی را افزایش دهد.
- absolute positional encoding:
  - شامل افزودن یا الحاق یک بردار به token embedding ورودی است که نشان دهنده موقعیت token در دنباله است. این موقعیت مستقل از موقعیت‌های دیگر در دنباله encode می‌شود. موقعیت‌های مطلق را نیز می‌توان در طی آموزش مدل یاد گرفت.
    - مزایا:
      - Simplicity: رمزگذاری برای هر موقعیت به طور مستقل ایجاد می‌شود و پیاده‌سازی را ساده می‌کند.
      - Effectiveness for small sequences: برای دنباله‌هایی که خیلی طولانی یا پیچیده نیستند، راه ساده‌ای برای اطلاع‌رسانی به مدل در مورد ترتیب توکن فراهم می‌کند.
    - معایب:
      - Limited context awareness: به طور ذاتی رابطه بین موقعیت‌های مختلف را نشان نمی‌دهند. آنها فقط برای مدل context مربوط به موقعیت‌های جدا از هم را به صورت تفکیک شده ارائه می‌دهند.
      - Scalability to Long Sequences: به خصوص با رمزگذاری سینوسی، مدیریت دنباله‌های بسیار طولانی می‌تواند مشکل‌ساز شود، زیرا ممکن است مدل به دلیل ماهیت تناوبی سینوسی‌ها، بین موقعیت‌های دور تمایز قائل نشود.
- Relative positional encoding:
  - فاصله یا رابطه بین token ها را در یک دنباله در نظر می‌گیرد. به جای encode کردن هر موقعیت به صورت مجزا، تفاوت در موقعیت‌های بین نشانه‌ها را encode می‌کند و بر نحوه ارتباط توکن‌ها با یکدیگر در چارچوب دنباله‌شان تمرکز می‌کند.

○ مزایا:

- Contextual relevance: با رمزگذاری روابط بین موقعیت ها، رمزگذاری های موقعیتی نسبی زمینه غنی تری را فراهم می کنند که می تواند به ویژه در کارهایی که نیاز به درک پویایی دنباله دارند، مانند مدل سازی زبان و تجزیه، مفید باشد.
- Better scalability for long sequences: رمزگذاری موقعیتی نسبی می تواند توالی های طولانی تر را به طور مؤثرتری مدیریت کند، زیرا مدل یاد می گیرد که بر تفاوت های موقعیتی که بیشترین ارتباط را برای کار دارند، بدون توجه به موقعیت مطلق در دنباله ورودی تمرکز کند.

○ معایب:

- Complexity: اجرای رمزگذاری موقعیتی نسبی به طور کلی پیچیده تر از رمزگذاری مطلق است. این نیاز به تغییراتی در مکانیزم self attention دارد تا تفاوت های موقعیت را در نظر بگیرد، که می تواند هم معماری مدل و هم پویایی آموزشی آن را پیچیده کند.
- Computational cost: رمزگذاری موقعیتی نسبی می تواند به دلیل محاسبات اضافی مورد نیاز در طول آموزش و inference، سرشار محاسباتی را افزایش دهد.

۲.

- Compatibility with Linear Self-Attention: یکی از مزیت های کلیدی RoPE سازگاری آن با مکانیسم های خطی Self-Attention است. رویکردهای سنتی برای relative positional encoding اغلب با چالش هایی در ادغام اطلاعات موقعیت نسبی به طور موثر با Self-Attention مواجه هستند. RoPE با ارائه روشی که به طور یکپارچه وابستگی های موقعیت نسبی را در فرمول توجه به خود ترکیب می کند، به این موضوع می پردازد، و توانایی مدل را برای درک روابط موقعیتی در دنباله ورودی افزایش می دهد.
- Decaying Inter-Token Dependency: RoPE یک ویژگی منحصر به فرد را معرفی می کند که در آن اطلاعات موقعیت نسبی با افزایش فاصله بین توکن ها کاهش می یابد. این کاهش در وابستگی بین token ها برای وظایف پردازش زبان طبیعی بسیار مهم است، زیرا به مدل کمک می کند تا روی وابستگی های محلی درون دنباله تمرکز بیشتری داشته باشد و در عین حال تأثیر token های دور را کاهش دهد. با ترکیب این ویژگی کاهشی، RoPE توانایی مدل را برای گرفتن اطلاعات متنی مرتبط در فواصل مختلف در توالی ورودی بهبود می بخشد و منجر به یادگیری و نمایش مؤثرتر می شود.
- Flexibility in Sequence Length: یکی دیگر از مزایای قابل توجه RoPE انعطاف پذیری آن در مدیریت توالی با طول های مختلف است. روش های رمزگذاری موقعیتی قبلی ممکن است با طول های دنباله ای متفاوت مشکل داشته باشند که بر عملکرد و سازگاری مدل با اندازه های ورودی مختلف تأثیر می گذارد. طراحی RoPE به آن اجازه می دهد تا به طور یکپارچه با توالی هایی با طول های مختلف سازگار شود و راه حلی قوی تر و همه کاره برای رمزگذاری اطلاعات موقعیت در مدل های transformer ارائه دهد. این انعطاف پذیری ظرفیت مدل را برای پردازش ورودی های با طول های مختلف بدون به خطر انداختن عملکرد یا کارایی افزایش می دهد.

۳. با ضرب ماتریس چرخش در بردار embedding، اندازه را ثابت را نگه میدارد و فقط زاویه بردار را تغییر میدهد. این به این معناست که پارامتر  $\theta$  که معادل relative distance، میان المان ها است تغییر میکند و در نتیجه با این عمل صرفاً relative position مربوط به عناصر بردار تغییر خواهد کرد.

۴.

۵. اصطلاح  $q_i K^T$  حاصل ضرب نقطه ای بردار  $q_i$  (query) با بردارهای key را نشان می دهد، که ارتباط بین query token و همه key token ها در دنباله ورودی را اندازه گیری می کند. عبارت  $[-(i-1), \dots, -2, -1, 0]$  یک bias مرتبط با موقعیت را بر



اساس فواصل بین query token و key token ها تعریف می کند. مقادیر منفی در دنباله نشان دهنده موقعیت نسبی key token ها با توجه به query token ها است. با افزودن عبارت bias خطی، مکانیسم attention تحت تأثیر موقعیت های نسبی token ها در دنباله ورودی قرار می گیرد. این bias اطلاعات مربوط به موقعیت را ارائه می دهد که مدل را برای توجه به موقعیت های مختلف بر اساس فاصله آنها از token پرس و جو راهنمایی می کند. پارامتر  $m$  نیز به عنوان head-specific slope، نقش scale کردن bias را دارد که قبل از زمان آموزش به صورت ثابت تنظیم میشود.

سوال پنجم:

$$Y = \text{Attention}(Q, K, V) = \text{Softmax}\left(Q \frac{K^T}{\sqrt{D}}\right)V \quad ۱.$$

در ساختار self-attention ماتریس های  $Q, K, V$ ، به ترتیب مربوط به query، key و value هستند که در آن ها هر سطر مربوط به بردار یک کلمه در دنباله ورودی است.  $D$  نیز ابعاد بردار های key است. این معادله میتواند به شکل ماتریسی نمایش داده شود که در آن هر دنباله ورودی از بردار های کلمات میتواند به طور پیوسته به بردار خروجی با ابعاد یکسان تبدیل شود.

۲. ماتریس های  $Q, K, V$  دارای ابعاد  $N \times D$  هستند ( $N$  طول دنباله ورودی و  $D$  ابعاد بردار کلمات). بنابراین ابعاد  $D \times N, K^T$  خواهد بود. در نتیجه ابعاد  $N \times N, QK^T$  خواهد شد. تابع softmax چون به صورت element wise انجام میشود تأثیری در اندازه خروجی ندارد. حاصل عبارت بعد از ضرب با ماتریس  $V$  دارای ابعاد  $N \times D$  خواهد شد. از آنجا که تعداد پارامتر های attention از مرتبه  $n^2$  است، در نتیجه تعداد کل پارامتر ها  $O(N^2 D^2) = N^2 D^2$  خواهد شد.

۳. در ماتریس self-attention، هر عنصر حاصل ضرب query یک کلمه در key کلمه دیگر است. اگر این ماتریس به صورت بلوکی در نظر گرفته شود که در آن هر بلوک معادل تعامل یک query و یک key باشد، میتوان دید که بسیاری از بلوک ها پارامتر هایشان را به اشتراک میگذارند. مخصوصاً زمانی که مقادیر وزن attention دو کلمه برای بقیه کلمات یکسان باشد، بلوک مربوط به تعامل آنها یکسان خواهد شد. همچنین به دلیل وجود softmax، تعداد زیادی از عناصر این ماتریس صفر خواهند شد.

۴. Positional encoding اطلاعات مربوط به ترتیب کلمات در دنباله ورودی را مشخص میکند. بنابراین اگر آن را حذف کنیم، ساختار self-attention فقط به معنای کلمات توجه میکند و مکان کلمات در جمله را نادیده میگیرد. در نتیجه خروجی ویژگی های مربوط به مکان کلمات را شامل نمیشود که منجر به همسانی با ترتیب تکراری ورودی می شود. بنابراین با تغییر ترتیب کلمات ورودی، ترتیب کلمات متناظر در خروجی نیز تغییر میکند که مطلوب نیست زیرا ترتیب کلمات در جمله اهمیت دارند.