

1.1

هر دو گزینه می توانند انتخاب های مناسبی باشند. به طور کلی LSTM از GRU قوی تر است و به نتایج برابر یا بهتری می رسد اما هزینه GRU ساده تر بودن آن است. GRU پارامترهای کمتری دارد و سریع تر است و هزینه محاسباتی کمتری دارد.

صمیمین می توان این طور در نظر گرفت که LSTM برای sequence های طولانی تر مفید است و long term dependency ها را بهتر حفظ می کند. بنابراین بسته به مسئله می توان بررسی کرد که آیا GRU به نتیجه دلخواه می رسد یا نه و اگر نتایج بهتری خواست و توان پردازشی آن را داشت از LSTM استفاده کند.

1.2

ضمیمه فرید کارت گرافیک را پیشنهاد نمی دهم. به احتمال بالا دلیل اصلی کندی محاسبات، ذات sequential معماری های LSTM و GRU است. با توجه به اینکه تسک این محقق sentiment analysis است، می دانیم که تمام توکن های هر sequence را به صورت یکجا در اختیار داریم اما با این دو معماری به صورت تدریجی آن ها را پردازش می کنیم.

به عنوان روش جایگزین می توان از attention و معماری transformer استفاده کرد که انگودر آن تمام توکن های ورودی را به صورت موازی پردازش می کند.

1.3.a

* ← ضرب element wise
@ ← ضرب ماتریسی

$$c^t = c^{t-1} * \Gamma_f^t + \Gamma_u^t * \tilde{c}^t$$

$$a^t = \tanh(c^t) * \Gamma_o^t$$

$$\hat{y}^t = \text{softmax}(w_y @ a^t + b_y)$$

$$\tilde{c}^t = \tanh \left(\overbrace{[w_{ca} \ w_{cn}]}^{w_c} @ \begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix} + b_c \right)$$

$K \times (K+d)$ $(K+d) \times 1$

$$\Gamma_o^t = \sigma \left([w_{oa} \ w_{on}] @ \begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix} + b_o \right)$$

$$\Gamma_u^t = \text{مشابه بالا}$$

$$\Gamma_f^t = \text{مشابه بالا}$$

مشتقات گزینش بردارهای a^t, c^t دارای طول k باشند و بردار x^t دارای طول d و بردار \hat{y}^t دارای طول l

1.3.6

$$\frac{\partial L}{\partial a^t} = \left[w_y^T \right]_{k \times l} @ \left[\hat{y}^t - y^t \right]_{l \times 1} + \left[da_{next} \right]_{k \times 1}$$

$$\frac{\partial L}{\partial c^t} = \left[\frac{\partial L}{\partial a^t} \right]_{k \times 1} * \left[\Gamma_o^+ \right]_{k \times 1} * \left[1 - \tanh^2(c^t) \right]_{k \times 1} + \left[dc_{next} \right]_{k \times 1}$$

$$\frac{\partial L}{\partial w_o} = \underbrace{\left(\left[\frac{\partial L}{\partial a^t} \right]_{k \times 1} * \left[\tanh(c^t) \right]_{k \times 1} * \left[\Gamma_o^+ * (1 - \Gamma_o^+) \right]_{k \times 1} \right)}_{\frac{\partial L}{\partial y_o^t}} @ \left[\begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix}^T \right]_{1 \times (k+d)}$$

$$\frac{\partial L}{\partial b_o} = \left[\frac{\partial L}{\partial y_o^t} \right]_{k \times 1}$$

$$\frac{\partial L}{\partial w_c} = \underbrace{\left(\left[\frac{\partial L}{\partial c^t} \right]_{k \times 1} * \left[\Gamma_u^+ \right]_{k \times 1} * \left[1 - \underbrace{\tanh^2(\rho \tilde{c}^t)}_{(\tilde{c}^t)^2} \right]_{k \times 1} \right)}_{\frac{\partial L}{\partial \rho \tilde{c}^t}} @ \left[\begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix}^T \right]_{1 \times (k+d)}$$

$$\frac{\partial L}{\partial b_c} = \left[\frac{\partial L}{\partial \rho \tilde{c}^t} \right]_{k \times 1}$$

$$\frac{\partial L}{\partial w_u} = \underbrace{\left(\left[\frac{\partial L}{\partial c^t} \right]_{k \times 1} * \left[\tilde{c}^t \right]_{k \times 1} * \left[\Gamma_u^+ * (1 - \Gamma_u^+) \right]_{k \times 1} \right)}_{\frac{\partial L}{\partial y_u^t}} @ \left[\begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix}^T \right]_{1 \times (k+d)}$$

$$\frac{\partial L}{\partial b_u} = \left[\frac{\partial L}{\partial y_u^t} \right]_{k \times 1}$$

1.3.6

$$\frac{\partial L}{\partial w_f} = \left(\underbrace{\left[\frac{\partial L}{\partial c^t} \right]_{k \times 1} * \left[c^{t-1} \right]_{k \times 1} * \left[\Gamma_f^t * (1 - \Gamma_f^t) \right]_{k \times 1}}_{\frac{\partial L}{\partial \gamma_f^t}} \right) @ \left[\begin{bmatrix} a^{t-1} \\ n^t \end{bmatrix}^T \right]_{1 \times (k+d)}$$

$$\frac{\partial L}{\partial b_f} = \left[\frac{\partial L}{\partial \gamma_f^t} \right]_{k \times 1}$$

$$\frac{\partial L}{\partial w_g} = \left[\hat{y}^t - y^t \right]_{\ell \times 1} @ \left[[a^t]^T \right]_{1 \times k}$$

$$\frac{\partial L}{\partial b_g} = \left[\hat{y}^t - y^t \right]_{\ell \times 1}$$

معمولاً برای upstream ها یک به عقب فرستاده می‌شود داریم :

$$\frac{\partial L}{\partial c^{t-1}} = \left[\frac{\partial L}{\partial c^t} \right]_{k \times 1} * \left[\Gamma_f^t \right]_{k \times 1}$$

$$\frac{\partial L}{\partial a^{t-1}} = \left[w_{oa}^T \right]_{k \times k} @ \left[\frac{\partial L}{\partial \gamma_o^t} \right]_{k \times 1} + \left[w_{ca}^T \right]_{k \times k} @ \left[\frac{\partial L}{\partial \rho \tilde{c}^t} \right]_{k \times 1} + \left[w_{ua}^T \right]_{k \times k} @ \left[\frac{\partial L}{\partial \gamma_u^t} \right]_{k \times 1} + \left[w_{fa}^T \right]_{k \times k} @ \left[\frac{\partial L}{\partial \gamma_f^t} \right]_{k \times 1}$$

2.1

$$l = \text{len}(\text{dict}) = 30,000$$

$$N = 2048 \quad \text{طول بیشترین دنباله ورودی}$$

$$h = 4 \quad \text{تعداد head ها}$$

$$d_k = \frac{768}{h} = 192$$

$$B = 32$$

$$d_{\text{model}} = 1024$$

همچنین فراموشی نکنید که در paper آن ذکر شده است، باید به این نکته در transformer

از residual connection استفاده شده است و همچنین اینکه ی دانیم layer norm ابعاد را تغییر

نمی دهد، بنابراین تمامی ورودی ها و خروجی ها دارای بُعد $(B \times N \times d_{\text{model}})$ هستند به غیر از

ورودی embedding در آنکودر و دیکودر که هنوز متن تبدیل به embedding ها نشده است

(پس از اینکه هر توکن با استفاده از dictionary به یک index تبدیل شود، این لایه آن را به

embedding تبدیل می کند پس متاثر بتواند ورودی لایه آیدینگ را $(1 \times N \times 1)$ در نظر گرفت)

index

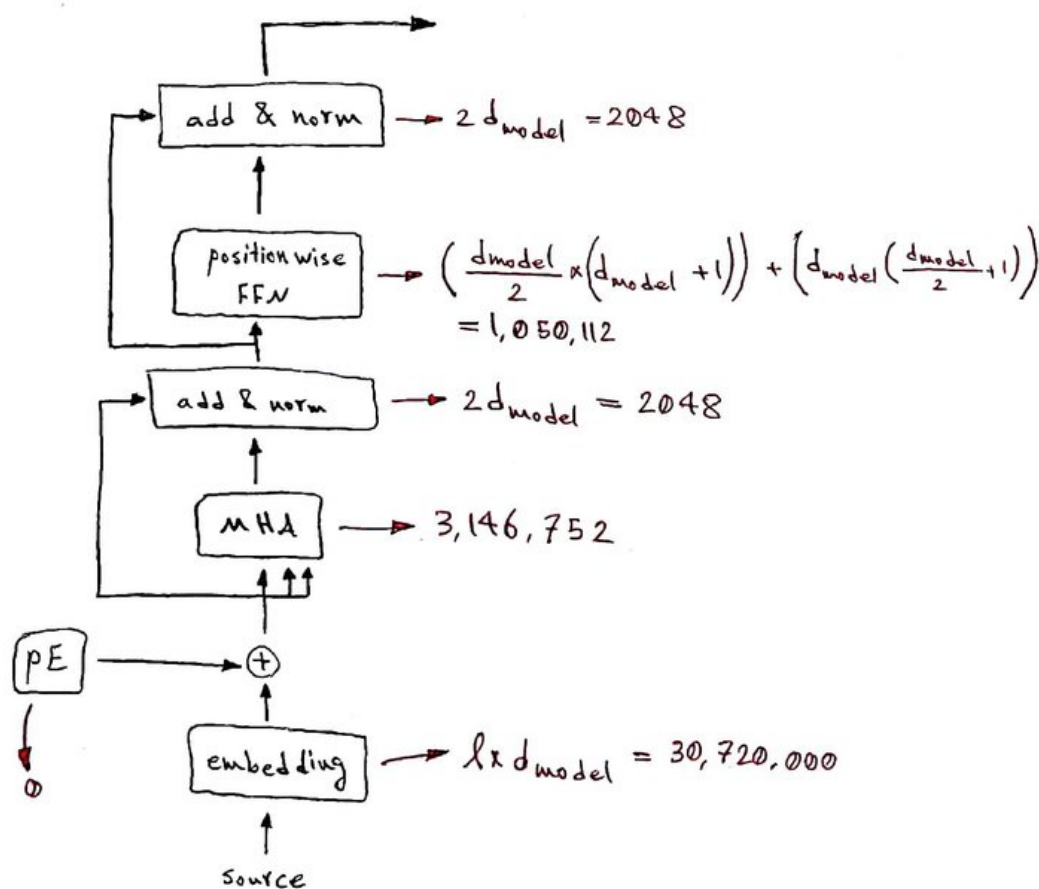
همچنین فراموشی نکنید که در این لایه نیز دارای بُعد $30,000$ بای به جای d_{model} است

بنابراین ابعاد خواسته شده در صورت سوال برابر است با:

$$(B \times N \times d_{\text{model}}) = (32 \times 2048 \times 1024)$$

2.2

encoder



multi head attention params:

$$h \left(d_{\text{model}} \times d_k \times 3 \right) \begin{pmatrix} Q \\ K \\ V \end{pmatrix}$$

$$+ d_{\text{model}} (d_k \times h) \quad \text{linear}$$

$$+ d_{\text{model}} \quad \text{linear bias}$$

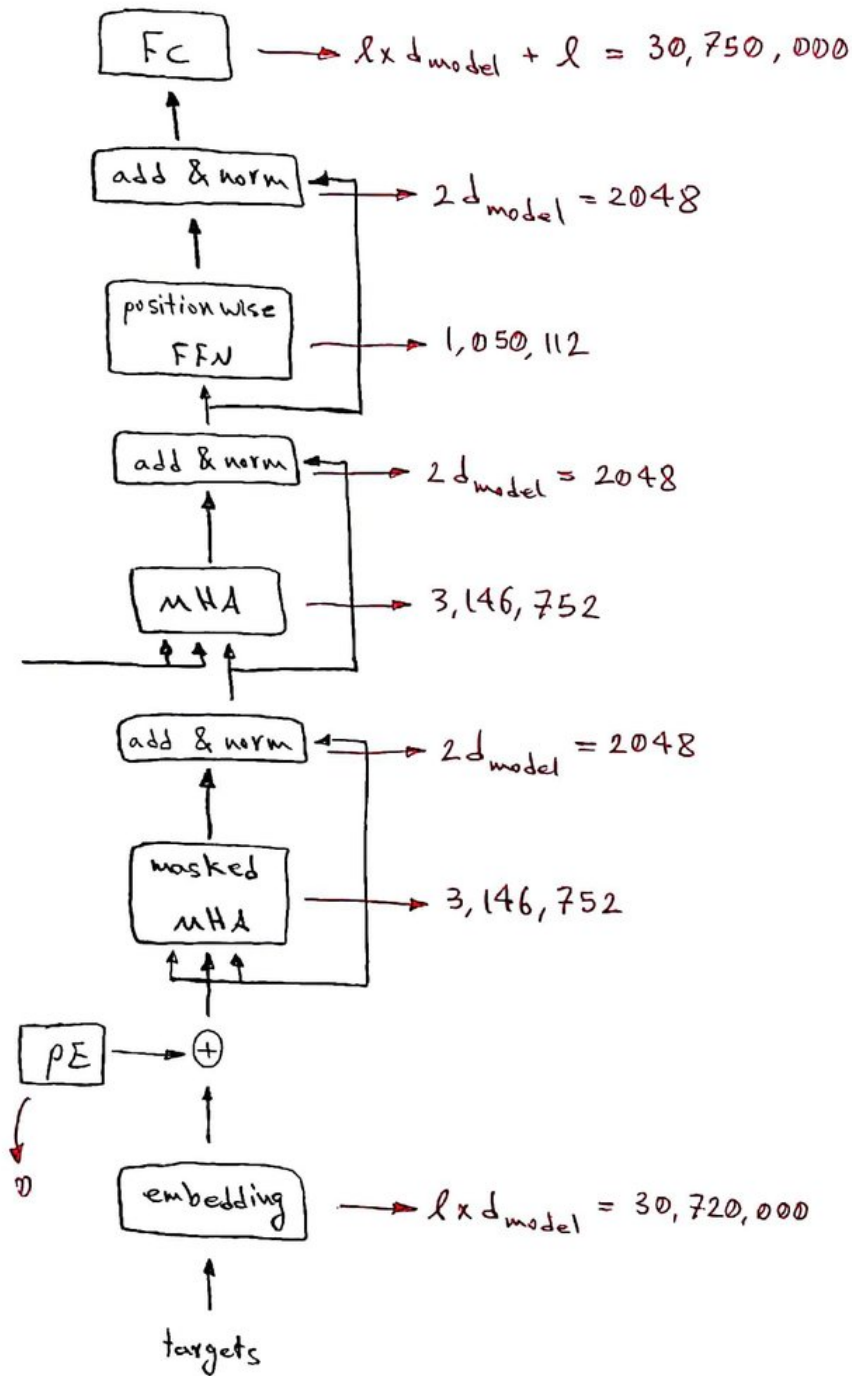
$$= 3,146,752$$

$$\text{encoder params} = 30,720,000 + 12 \left(4,200,960 \right) = 81,131,520$$

2.2

: f/b encoder كسب

decoder



$$\text{decoder params} = 30,720,000 + 30,750,000 + 8(7,349,760) = 120,268,080$$

$$\text{encoder params} + \text{decoder params} = 201,399,600$$

اگر تمام بردارهای دوری به هم عمود باشند داریم:

3.1 $i \neq j \quad a_{ij} = \frac{1}{\underbrace{1 + \dots + 1}_{N-1} + \|x_i\|^2} \approx 0$

$a_{ii} = \frac{\|x_i\|^2}{1 + \dots + 1 + \|x_i\|^2} \approx 1$

$A \approx I_N \rightarrow y_n = \sum_{m=1}^N a_{nm} x_m$

$= a_{nn} x_n \approx x_n$

بنابراین بردارهای خروجی (y_n) تقریباً همان بردارهای دوری (x_n) می‌شوند

3.2

ایمان‌های بردارهای a, b, c با a_i, b_i, c_i نشان می‌دهیم

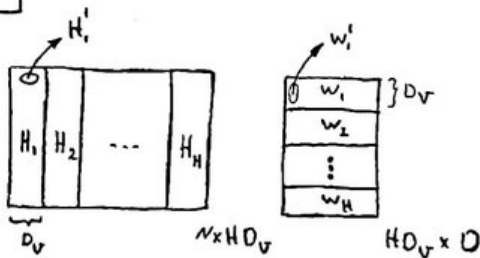
$a_i, b_i \sim \mathcal{N}(0, 1)$

$E[(a^T b)^2] = E[(a_1 b_1 + \dots + a_D b_D)^2] \stackrel{\text{I}}{=} \text{Var}[a_1 b_1 + \dots + a_D b_D] \stackrel{\text{استقلال}}{=} \underbrace{\text{Var}[a_1 b_1]}_1 + \dots + \underbrace{\text{Var}[a_D b_D]}_1 = D \stackrel{\text{II}}{=}$

$E[a_1 b_1 + \dots + a_D b_D] = \sum_{i=1}^D E[a_i b_i] \stackrel{a_i \perp b_i}{=} \sum_{i=1}^D \underbrace{E[a_i]}_0 \underbrace{E[b_i]}_0 = 0 \stackrel{\text{I}}{=}$

$\text{Var}[a_i b_i] = E[a_i^2 b_i^2] - \underbrace{E[a_i b_i]^2}_0 \stackrel{a_i \perp b_i}{=} E[a_i^2] E[b_i^2] = \text{Var}[a_i] \text{Var}[b_i] = 1 \stackrel{\text{II}}{=}$

3.3



⊗ ماتریس $W^{(0)}$ در صورت سوال با W می‌تایم تا نزدیک ما خلوت تر شوند!

سطر نام ماتریس H_j با H_j^T می‌تایم که دارای بُعد $1 \times D_V$ است

ستون نام ماتریس w_j با w_j^T می‌تایم که دارای بُعد $D_V \times 1$ است

ضرب 2 ماتریس
بالا به صورت بلوک

$$= \begin{bmatrix} \sum_{i=1}^H H_1^T w_1^T & \dots & \sum_{i=1}^H H_1^T w_D^T \\ \vdots & & \vdots \\ \sum_{i=1}^H H_N^T w_1^T & \dots & \sum_{i=1}^H H_N^T w_D^T \end{bmatrix}_{N \times D} = \sum_{i=1}^H \begin{bmatrix} H_1^T w_i^T & \dots & H_1^T w_i^D \\ \vdots & & \vdots \\ H_N^T w_i^T & \dots & H_N^T w_i^D \end{bmatrix}_{N \times D}$$

$$= \underbrace{\begin{bmatrix} H_1^T \\ \vdots \\ H_N^T \end{bmatrix}}_{N \times D_V} \underbrace{\begin{bmatrix} w_1^T & \dots & w_D^T \end{bmatrix}}_{D_V \times D}$$

$$= H_i^T W_i$$

3.3

$$Y(x) = \sum_{i=1}^H H_i w_i = \sum_{i=1}^H \text{softmax} \left(\frac{Q_i K_i^T}{\sqrt{D_{K_i}}} \right) \underbrace{w_i^{(v)}}_{w^{(i)}} w_i$$

4.1

learnable vs fix

در حالت learnable جاسازی موقعیت دارای پارامترهایی است که مانند سایر پارامترهای مدل learn می‌شوند و بنابراین این روش انعطاف بیشتری در برابر داده‌های مختلف و شرایط مختلف دارد اما باید عوالممان باشیم که پارامتر بیشتر واحدی به معنای d_{vec} بیشتر است و بنابراین باید مقدار داده‌ها نیز به اندازه کافی باشد وگرنه گرفتار overfitting می‌شویم.

در حالت fix جاسازی موقعیت فاقد پارامتر بوده و یک تابع fix است که از روی فرمولی از پیش تعریف شده، فردی تولید می‌کند. به عنوان مثال در معماری ترنسفورمر از جاسازی موقعیت ثابت استفاده شده است. این روش انعطاف کمتری نسبت به شرایط و مسائل مختلف دارد زیرا learn نمی‌شود اما چون ثابت و از پیش تعریف شده است، مدل را ساده‌تر می‌کند و سریع‌تر است. البته باید در نظر داشت که توابع ثابت نیز می‌توانند inductive bias های مناسبی (مطابق کاربرد) در مدل تزریق کنند که کمک کننده باشند.

absolute vs relative

در حالت absolute برای جاسازی موقعیت به دست می‌آید تنها بر اساس موقعیت مکانی یا همان ایندکس توکن در sequence است. با اینکه به نظر می‌رسد همین موقعیت مکانی در توکن به تنهایی کافی باشد اما مشاهده می‌شود که در این روش موقعیت نسبی توکن‌ها نسبت به هم، آن طور که باید، capture نمی‌شود. به عنوان مثال انتظار داریم که فاصله positional emb. دو توکن که از هم دورتر هستند بیشتر از فاصله positional emb. دو توکن که به هم نزدیک هستند باشد. اما این حالت همیشه رخ نمی‌دهد (به خصوص در حالت learnable که inductive bias ای برای این نقشه نداریم) در حالت relative به نسبت بین توکن‌ها بیشتر پرداخته می‌شود تا مشکل ذکر شده حل شود. به عنوان مثال اگر یک sequence چهار کلمه‌ای داشته باشیم و یک کلمه به آن امانه کنیم، positional emb. در 4 کلمه قبلی عوض می‌شود. در این روش چون تمام pair های sequence نسبت به هم سنجیده می‌شوند محاسبات سنگین‌تری داریم. هم در زمان آموزش و هم در زمان تست فرایند کمتری خواهیم داشت. در پیاده‌سازی حالت relative نیز برخلاف absolute که با embedding اولیه جمع می‌شود، در attention score ما اعمال می‌شود.

4.2

در روش rope سبب شده است مزیت های هر دو روش absolute , relative را ترکیب کنیم . هم از ایندکس توکن ها در sequence تأثیر بپذیریم و هم موقعیت ها به صورت نسبی capture شوند و هم محاسبات efficient ای داشته باشیم !
 اعمال این positional emb. از طریق ضرب یک ماتریس rotation در embedding ها صورت می پذیرد که زاویه چرخش $m\theta_1$ است . m بخشی است که مشابه حالت absolute است و همان ایندکس توکن ها (صرف نظر از نسبت ها) است .
 از طریقی چون کلمات پایانی نسبت به کلمات ابتدایی m بزرگتری دارند ، زاویه چرخش بیشتری هم دارند و بنابراین نسبت ها نیز تأثیرگذار خواهند بود .
 همچنین محاسبات نیز دارای پیچیدگی کمتری است و به خصوص در زمان تست سرچایی نداریم . چنی هر 3 مشکل ذکر شده تا حد مناسبی بهبود یافته .

4.3

اندازه بردار معادل embedding اولیه است که هنوز positional emb. به آن افزوده نشده است
 زاویه بردار معادل embedding ای است که positional emb. به آن افزوده شده است
 بنابراین روش rope به جای جمع کردن بردار positional emb. با ایندکس قبلی ، اطلاعات را در دو بُعد اندازه و زاویه با هم ترکیب می کند به صورتی که ماتریس rotation اندازه را تغییر نمی دهد

4.4

روش سینوسی fix است و پارامتری ندارد اما روش rope یک روش learnable است .
 به طور کلی طبق مزیت های ذکر شده روش rope اطلاعات relative و absolute موقعیت های مکانی را بهتر capture می کند و به دلیل learnable بودن انعطاف پذیری بیشتری هم دارد . همچنین روش ساده سازی این دو روش نیز فرق دارد ، در حالت سینوسی اطلاعات positional emb. با ایندکس اولیه جمع می شود اما در روش rope ماتریس های rotation در K, Q ضرب می شوند

برای حالت سینوسی داریم :

$$PE(pos, z_i) = \sin(pos / 10000^{\frac{2i}{d_{model}}})$$

$$PE(pos, z_{i+1}) = \cos(pos / 10000^{\frac{2i}{d_{model}}})$$

برای مثال d_{model} را برابر 4 می گیریم :

$$\begin{bmatrix} \sin 0 \\ \cos 0 \\ \sin 0 \\ \cos 0 \end{bmatrix}$$

"دانشگاه"

$$\begin{bmatrix} \sin 1 \\ \cos 1 \\ \sin \frac{1}{10000^{\frac{2}{4}}} \\ \cos \frac{1}{10000^{\frac{2}{4}}} \end{bmatrix}$$

"صنعتی"

$$\begin{bmatrix} \sin 2 \\ \cos 2 \\ \sin \frac{2}{10000^{\frac{2}{4}}} \\ \cos \frac{2}{10000^{\frac{2}{4}}} \end{bmatrix}$$

"شریف"

⊗ ایندکس های به دست آمده با ایندکس های اولیه جمع می شوند به صورت element wise

4.4 ^{اول}

$$\theta_i = 10000 \frac{-2(i-1)}{d}$$

$$\rightarrow \theta_1 = 1$$

$$\theta_2 = 10000 \frac{-2}{4} = \frac{1}{100}$$

برای حالت $Rope$ داریم: (نرخ می کشیم $d=4$)

⊕ ماتریس های به دست آمده از سمت راست در ماتریس های Q ، K ضرب ماتریسی می شوند

$$\begin{bmatrix} \cos \theta & -\sin \theta & 0 & 0 \\ \sin \theta & \cos \theta & 0 & 0 \\ 0 & 0 & \cos \theta & -\sin \theta \\ 0 & 0 & \sin \theta & \cos \theta \end{bmatrix} = I$$

"کوانتا"

$$\begin{bmatrix} \cos 1 & -\sin 1 & 0 & 0 \\ \sin 1 & \cos 1 & 0 & 0 \\ 0 & 0 & \cos \frac{1}{100} & -\sin \frac{1}{100} \\ 0 & 0 & \sin \frac{1}{100} & \cos \frac{1}{100} \end{bmatrix}$$

"سین"

$$\begin{bmatrix} \cos 2 & -\sin 2 & 0 & 0 \\ \sin 2 & \cos 2 & 0 & 0 \\ 0 & 0 & \cos \frac{2}{100} & -\sin \frac{2}{100} \\ 0 & 0 & \sin \frac{2}{100} & \cos \frac{2}{100} \end{bmatrix}$$

"سین"

4.5

در نرم ماتریسی برای حالت auto regressive داریم :

$$\text{softmax} \left(\begin{bmatrix} q_1 \cdot k_1 & & & \\ q_2 \cdot k_1 & q_2 \cdot k_2 & & \\ & & \ddots & \\ q_N \cdot k_1 & q_N \cdot k_2 & \dots & q_N \cdot k_N \end{bmatrix}_{N \times N} + m \begin{bmatrix} 0 & & & \\ -1 & & 0 & \\ \vdots & & & \\ -(N-1) & -(N-2) & \dots & 0 \end{bmatrix}_{N \times N} \right)$$

روش عملکرد AliBi به این صورت است که یک ماتریس QK^T افاده می کند که بر اساس نامده هر جفت توکن از یکدیگر به شباهت آن ها پالتی می دهد . یعنی اگر 2 کلمه شباهت زیادی داشته باشند ، ماتریس اول مقدار ضرب داخلی بیشتری خواهد داشت اما اگر آن دو کلمه از هم دور باشند بخش ماتریس به اندازه نامده آن دو کلمه از مقدار شباهت آن ها کم می کند که باعث می شود attention score نسبت به هم داشته باشند .

به این صورت توسط یک روش غیر learnable توانستیم نسبت بین توکن ها را نیز در معادله دخیل کنیم .

توسط m می توان میزان تأثیر نسبت را کنترل کرد . هر چه m بیشتر باشد ، نزدیکی بودن کلمات به هم مهم تر از شباهت امپدینگ آنها می شود .

\otimes m به ازای هر head جدا تعیین می شود که قابلیت کنترل بیشتری می دهد

5.1 طبق توضیحات ارائه شده در کد بالا، X را به صورت یک بردار $N \times 1$ در نظر می‌گیریم و برای عبارت داخل softmax داریم:

$$\frac{QK^T}{\sqrt{D}} = \frac{1}{\sqrt{D}} \times \underbrace{\begin{bmatrix} w_Q \\ \vdots \\ w_Q \end{bmatrix}}_{\substack{N \times D \\ N \times 1}} @ \underbrace{\begin{bmatrix} X \\ \vdots \\ X \end{bmatrix}}_{\substack{N \times D \\ 1 \times N}} @ \underbrace{\begin{bmatrix} X^T \\ \vdots \\ X^T \end{bmatrix}}_{\substack{1 \times D \\ N \times 1}} @ \underbrace{\begin{bmatrix} w_K^T \\ \vdots \\ w_K^T \end{bmatrix}}_{\substack{N \times D \\ N \times 1}}$$

حال XX^T را ورودی در نظر می‌گیریم. با توجه به اینکه ماتریس w_K^T مربعی است می‌دانیم که می‌توانیم برای آن، ماتریس $[w'_K]$ را یافت که ضرب کردن w'_K از چپ در عبارت فوق معادل ضرب کردن w_K^T از راست در آن عبارت شود. یعنی:

$$\frac{QK^T}{\sqrt{D}} = \frac{1}{\sqrt{D}} \underbrace{w'_K w_Q}_W \underbrace{XX^T}_{\text{ورودی}}$$

بنابراین نشان دادیم که عبارت داخل softmax معادل یک ضرب ماتریسی است.

ممکن است softmax را نیز می‌توان معادل activation function ای در نظر گرفت که روی این لایه خطی اعمال می‌شود.

$$Y = \text{softmax} \left(\underbrace{\begin{bmatrix} w \\ \vdots \\ w \end{bmatrix}}_{\substack{N \times D \\ N \times 1}} @ \underbrace{\begin{bmatrix} XX^T \\ \vdots \\ XX^T \end{bmatrix}}_{\substack{N \times D \\ N \times 1}} \right) @ \underbrace{\begin{bmatrix} v \\ \vdots \\ v \end{bmatrix}}_{\substack{N \times D \\ 1 \times 1}}$$

بعد از محاسبه ضربی نیز معادل بعد ورودی است زیرا:

پس Y نیز مانند X دارای بُعد $N \times 1$ است.

5.2 طبق بحث قبلی ماتریس w دارای بُعد $N \times N$ است و بنابراین تعداد پارامترها $O(N^2 D^2)$ می‌شود

5.4 تفاوت multi head attention در حالتی که یک head داشته باشیم در طول اینترنل کلمات (در حالت میانی) است و بنابراین در این سوال که با تغییر ترتیب کلمات کار دارد، تفاوتی ایجاد نمی‌کند. در حالت کلی طبق فرمول attention داریم:

$$QK^T = \begin{bmatrix} x_1^T w_Q \\ \vdots \\ x_N^T w_Q \end{bmatrix}_{N \times D_K} @ \begin{bmatrix} | & & | \\ w_K^T x_1 & \dots & w_K^T x_N \\ | & & | \end{bmatrix}_{D_K \times N}$$

واضح است که ضرب ماتریسی بالا نسبت به جایگذاری کردن x_i ها (ترتیب کلمات) equivariant است زیرا سیاهه هر کلمه نسبت به تمام کلمات سنجیده می‌شود و به همان ترتیب کلمات ورودی در ماتریس $N \times N$ حاصل شده قرار می‌گیرد. حال اگر ترتیب ورودی تغییر کند، طبق ترتیب جدید، تمام سیاهه ها به دست آمده و در ماتریس $N \times N$ حاصل شده تکرار می‌گردد.

اعمال softmax که اثری در این قضیه ندارد و ضرب در ماتریس $\begin{bmatrix} x_1^T w_V \\ \vdots \\ x_N^T w_V \end{bmatrix}_{N \times D_K}$ نیز دنیا طبق استدلال ذکر شده نسبت به جایگذاری کردن x_i ها equivariant است.

5.4 ادانه بنا بر این همان طور که ذکر شد multi head کردن سپس concat کردن در بُعد طول امیدواریم هر کلمه است نه ترتیب
 کلمات و هیچ اثری در این جایابی ندارد و تمام ضرب های ماتریسی صورت گرفته نیز نسبت به جایابی ورودی ها equivariant
 هستند. بنا بر این هر ترتیبی که در X داشته باشیم، همان ترتیب سطر ها را نیز در Y داریم و ترتیب توالی های sequence عملاً لحاظ نمی شود.
 برای همین است که از positional encoding استفاده می شود.

❗ در برخی سوال ها با علی بابایک همفکری داشتیم