

به نام خدا

علی قاسم زاده 401106339

به نام خدا

علی قاسم زاده 401106339

۱- الف) ابتدا تعداد Positive pair ها را نسبت می آوریم

هر کلام از کلمات اول و دوم به ~~۲ عدد~~ ۲ عدد جفت مثبت تشکیل می دهد
هر کلام از کلمات دوم و سوم ۲ عدد جفت مثبت تشکیل می دهد. برابر سایرین هر ۴ تا
جفت مثبت داریم:

$$\text{تعداد جفت ها / مثبت} \rightarrow 2 \times 2 + 2 \times 3 + 6 \times 4 = 34$$

$$\text{حالا برابر هر سمبل مثبت ۵ سمبل منفی داریم} \\ \text{جفت} \quad (5+1) \times 34 = 204$$

۳۰۰. اگر target word مان در جایگاه i قرار گیرد، مان را L در نظر می گیریم
Corups لای

$$\text{تعداد کلمات ورودی} = \min(i-1, 4) + \min(L-i+4)$$

d-model برابر است با ابعاد input embedding h هر تعداد کلمات attention هاست.

$$d_k = d_v = \frac{d_{\text{model}}}{h}$$

ابعاد ماتریس ها / Projection مان به صورت زیر است:

$$\textcircled{a} \quad w_i^Q \in \mathbb{R}^{d_{\text{model}}} \quad w_i^K \in \mathbb{R}^{d_{\text{model}}} \quad w_i^V \in \mathbb{R}^{d_{\text{model}}}$$

این تقسیم باعث می شود که رویش کوچکتر از مقدار برداری اصلی، مما سبت مان انجام شود
بعث می شود مما سبت ماتریسی مان دیدار از سی معاری بصورت انجام شود

ب. درست: در معاری اصلی transformer از fixed sinusoidal positional encodings استفاده

~~بیشتر~~ از پیش هر تحسین شده است و یادگیر نمی شود (train نمی شود)

در تابع این encode کردن برابر این به کار می رود که ترتیب کلمات اطلاعات حاصل از

را نیز نگه داریم. چون مکانیزم attention به ترتیب token ها اهمیت نمی دهد.

ولی Bert، positional embedding است قابل یادگیر است.

② embedding استی را هم باید یادگیر باشد trainable مناره
← آموزش مدل بهتر انجام می شود.

• غلط است، همواره عبارت صورت سوال درست نیست زیرا:

هرچقدر تعداد head ها بیشتر کنیم، ویژگی ها را بیشتر

استخراج می شود ولی باعث افزایش هزینه/مسابقات و اتفاقاتی مثل

overfitting می شود، همچنین ما قله/بیشتر نیاز خواهیم داشت و

resource بیشتر برابر مسابقات نیاز خواهیم داشت.

$$\begin{aligned}
 - \frac{\partial \log P(w_o | w_t)}{\partial v_{w_o}} &= - \frac{\partial \log P(w_o | w_t)}{\partial P(w_o | w_t)} \cdot \frac{\partial P(w_o | w_t)}{\partial v_{w_o}} = \\
 &= - \frac{1}{P(w_o | w_t)} \cdot \frac{\partial \left(\frac{\exp(u_{w_t}^T v_{w_o})}{\sum_{k \in V} \exp(u_{w_t}^T v_k)} \right)}{\partial v_{w_o}} = A \\
 &= - \frac{1}{P(w_o | w_t)} \cdot \frac{u_{w_t} \exp(u_{w_t}^T v_{w_o}) A - u_{w_t} \exp(u_{w_t}^T v_{w_o})^2}{A^2} \\
 &= - \frac{1}{P(w_o | w_t)} \cdot u_{w_t} \left(P(w_o | w_t) - P(w_o | w_t)^2 \right) \\
 &= u_{w_t} (P(w_o | w_t) - 1)
 \end{aligned}$$

ب)

۱-۲: مدل‌های مختلف حساب‌های خیلی نزدیک‌تر هستند و وابسته هستند به این ساختار
 برای داده ساختارهایی که داده‌ها همگی دارند خوب عمل می‌کند، ارتباط‌ها/طولانی
 مدت را هم می‌توانید می‌گیرید. می‌سازد هر کس و آسان‌تر خواهد بود چرا که
 تعداد neg pair, pos pair کمتر خواهند بود.

۳-۴: مدل‌ها به تعداد همسایه‌های نزدیک بستگی دارد و مدل‌ها می‌تواند ارتباطات بین
 کلمات اطراف را خوب یاد بگیرد، معمولاً هم برای داده‌ها/عوی خوب عمل
 می‌کند. هزینه‌های سببش هم این است، اما این است می‌تواند متوسط
 در نظر گرفته می‌شود.

۵-۶: مدل‌ها به تعداد همسایه‌ها/زیادی و روابط طولانی بین کلمات توجه می‌کند
 و دنباله‌های طولانی‌تر را در نظر می‌گیرد، از طرفی زمان زیادی هر طول می‌کشد تا مدل بتواند
 شود و تعداد sample بیشتر بهر یادگیر نیاز خواهد داشت.

3.1. اگر دنباله ی ورودی بردارهای کلمه به شکل $X = \{x_1, x_2, \dots, x_n\}$ باشد آنگاه داریم :

$$\text{Attention}(Q, K, V) = \text{softmax}(QK^T/\sqrt{d_k})V$$

حال می‌خواهیم عملیات را به صورت یک لایه ی fc در نظر بگیریم :

بردار های x_i را در کنار هم قرار می دهیم تا ماتریس Y را بسازیم، داریم که :

$$Q=YW_Q, \quad K=YW_K, \quad V=YW_V$$

ماتریس A را به این صورت می سازیم که $A = YY^T$:

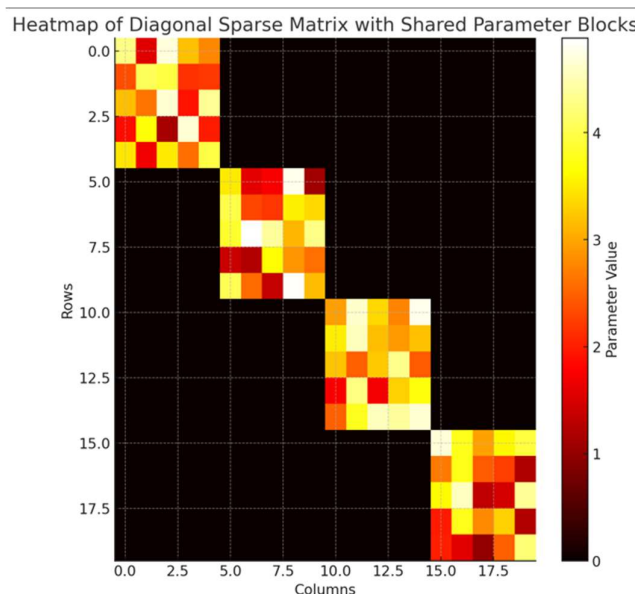
این ضرب هم روابط بین بردار های ورودی را نشان می دهد. حالا برای مپ کردن ورودی به یک بردار خروجی، یک ضرب بین این ماتریس و ماتریس وزن های W انجام می دهیم که W قابل یادگیری است.

$$\text{Softmax}(AW)$$

این عبارت هم معادل یک FC با تابع $\text{activation} = \text{softmax}$ است.

2.3. در شبکه قبلی داریم که W ماتریسی $N*N$ است و در نتیجه $O(N^2)$ تا پارامتر دارد. حالا می خواهیم $N*N$ را به $D*D$ ببریم پس به $O(N^2D^2)$ پارامتر نیاز داریم.

3.3. در مکانیزم attention به هر کلمه در دنباله ی ورودی توجه می شود ولی میزان توجه مان به یک اندازه نیست و وزن های متفاوتی داریم همچنین چون از softmax استفاده می کنیم توزیع توجه مان پراکنده خواهد شد و تعداد زیادی از وزن ها نزدیک صفر خواهند ماند در نتیجه ماتریسی sparse ایجاد می شود که به کلمات مهمتر بیشتر وزن برای توجه می دهد.



4. مکانیزم attention نسبت به ترتیب ورودی ها اهمیتی نمی دهد، یعنی اگر جابه جا کنیم ترتیب ورودی ها را باز هم همان خروجی را می دهد برای همین نیاز داریم تا اطلاعات مکانی کلمات را در جایی ذخیره کنیم برای همین از positional encoding استفاده می کنیم.

اگر ماتریس X را با ابعاد $N \times D$ در نظر بگیریم، حالا جایگشت p را در نظر می گیریم که ترتیب سطر ها را عوض می کند. خواهیم داشت که $X_p(i)$ حالا داریم که :

$$Q' = YW_Q, K' = YW_K, V' = YW_V$$

$$Q_p(i) = X_p(i)W_Q, K = X_p(i)W_K, V = X_p(i)W_V$$

حالا باید داشته باشیم که :

$$(QK^T)_{ij} = (Q'K'^T)_{p(i)p(j)}$$

حالا softmax را اعمال می کنیم، خواهیم داشت که :

$$\text{Softmax}(Q'K'^T/\sqrt{d_k})_{p(i)p(j)} = \text{Softmax}(QK^T/\sqrt{d_k})_{ij}$$

حالا وزن های بدست آمده برای هر کلمه در این softmax دقیقاً متناظر با همان وزن هایی هستند که بدون جایگشت داشتند. در نتیجه اگر جایگشت بدهیم تغییر در وزن خروجی i ام ایجاد نمی شود

$$O_i = O_p(i)$$

در نتیجه حکم اثبات می شود .

encoder inp, $[32, 2048, 1024]$	decoder inp, $[32, 2048, 1024]$ ① - ٢
Per-head self att, $[32, 2048, 192]$	Per-head self att, $[32, 2048, 192]$
all heads, $[32, 2048, 1024]$	all heads, $[32, 2048, 1024]$
Forward 1, $[32, 2048, 1024]$	Per-head cross-att, $[32, 2048, 192]$
Forward 2, $[32, 2048, 1024]$	all cross-att, $[32, 2048, 1024]$
Encoder, $[32, 2048, 1024]$ ^{out}	Forward 1, $[32, 2048, 1024]$
	Forward 2, $[32, 2048, 1024]$
	out, $[32, 2048, 20000]$

embedding vector

$$d_k = d_v = \frac{\sqrt{48}}{4} = 192$$

② - ٢
self attention

$$\text{Param}(W_Q) = \text{Param}(W_K) = \text{Param}(W_V) = 32 \times 2048 \times 192 = 125829376$$

$$\text{Param}(W_O) = 48 \times 1024 = 49152$$

Multi-head attention / تفاعل بالمشغول

$$M \xrightarrow{\text{self}} 32 \times (32 \times 2048 \times 192) + 48 \times 1024 = 3145728$$

Feed Forward / تفاعل بالمشغول

$$F_1 \rightarrow \text{Param}(W_1) = 1024 \times 1024 + 1024 = 1048576$$

$$F_2 \rightarrow \text{Param}(W_2) = 1024 \times 1024 + 1024 = 1048576$$

تفاعل بالمشغول / تفاعل بالمشغول

$$B = M + F_1 + F_2 = 49152$$

encoder / تفاعل بالمشغول

$$E = A \times B = 32 \times 2048 \times 20000$$

حالا تعداد پارامترها / decoder را حساب می‌کنیم.
 دیکودرمان یک لایه cross att نیز دارد که تعداد پارامترها آن برابر
 است با تعداد پارامترها / self-att

$$12 \times (41958480 + 3145728) = 18098812$$

تعداد پارامترها / embedding برابر لایه آخر:

$$1024 \times 30000 + 30000 = 30750000$$

تعداد کل پارامترها عبارت است از:

$$33524720 + 18098812 + 30750000 = 82373532$$

3. "جمله" علی به مربع می‌رود "را در" نقل می‌کنیم، خواسته درست که:

ابتدا هر کلمه در این جمله به برطری با 1024 می‌شود که اطلاعات آن
 مثل اطلاعات معنایی، دستوری و ... را encode می‌کند.

در واقع ~~هر کلمه~~ جمله ما به تعداد token تقسیم می‌شود و سپس

هر توکن به یک بردار 1024 تایی می‌شود

در مرحله tokenize کردن، اطلاعات positional

نیز به بردارها اضافه می‌شوند تا اطلاعات ترتیب کلمات حفظ شود

در مرحله encoder

برای هر توکن ورودی ماتریس ها

$K \rightarrow \text{key}$

$V \rightarrow \text{value}$

$Q \rightarrow \text{query}$

حاسب می شوند.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

attention از رابطه بالا استفاده می شود. برای هر Head مقدار Attention محاسب می شود

و در نهایت تمامی مقادیر ادغام می شوند و در ماتریس W_0 ضرب می شوند

$$W_0 \in \mathbb{R}^{768 \times 1024}$$

نتایج هر از دو لایه Forward عبور می کنند (با تابع فعال سازی σ)

$$\text{FFN}(x) = \sigma(xw_1 + b_1)w_2 + b_2$$

خروجی encoder مثال خروجی دومین Forward است

پس دیگر در self attention را مشابه Encoder برای توکن ها حساب

می کنند پس برای ارتباط بین ورودی دیگر در ~~(دو لایه)~~ و خروجی انکودر از cross-attention استفاده می شود

$$\text{Attn}(Q, K, V) = \left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

نتایج در نهایت از دو لایه Forward عبور می کنند و به بردار خروجی می گیریم.

پس با استفاده از تابع softmax احتمال ها را روی کلمات حساب می شود

در نهایت کلمات به ~~این~~ صورت پیش می روند که

کلمات / اول پیش بینی می شود پس کلمات / بعد به شرط آن
د الی آخر. (هر کس به شرط تمام قبلی ها است)
اینقدر کلمات تولید می شود تا به پایان جمله و حد اکثر طول کلمات برسیم.

5. الف) مدل bert از NSP در زمان آموزش خود استفاده می کند، این کمک می کند تا ارتباط بین جملات متوالی را بتواند بهتر درک کند به صورت عمل می کند که با تعدادی pair sentence مدلمان train می شود و بررسی میکند که جمله ی دوم چقدر به اولی وابسته است یا نه این کمک می کند تا ارتباط بین جملات و ترتیب آنها را بیاموزد.

برای تسک های پایین دستی که نیاز به درک معنای کلی از متن دارند نیازی نداریم که ارتباط بین جملات را بررسی کنیم و با حذف nsp بیشتر می توانیم روی درک مفهوم کلی جمله تمرکز کنیم و در نتیجه عملکردمان بهبود می یابد.

ب) **تنظیم دقیق (Fine-tuning) مدل هایی مانند BERT** برای تسک های پایین دستی (downstream tasks)، زمانی که نیاز به درک کلی جمله و معنای آن در سطحی وسیع تر داریم، چالش های خاص خود را به همراه دارد. این مسئله به ویژه در شرایطی مطرح می شود که مدل فقط برای تسک Masked Language Modeling (MLM) آموزش دیده باشد و فاقد قابلیت های پیش فرض برای درک روابط کلی جمله باشد (مانند مدل RoBERTa که NSP را حذف کرده است).

استفاده از لایه های Pooling برای استخراج نمای کلی جمله:

- **ایده اصلی:** در این روش، به جای استفاده از بردارهای مرتبط با هر توکن به صورت جداگانه، یک نماینده ثابت (بردار معنادار) برای کل جمله استخراج می شود. این نمایه می تواند خلاصه ای از اطلاعات موجود در جمله باشد و نمایانگر معنای کلی آن باشد.

- **روش های رایج در لایه های Pooling:**

1. **Mean Pooling:** در این روش، میانگین تمام بردارهای خروجی کلمات جمله گرفته می شود. این میانگین می تواند به عنوان نمای کلی جمله استفاده شود.

2. **Max Pooling:** در این روش، بیشترین مقدار در هر بُعد از بردارهای خروجی کلمات انتخاب می شود. این روش ویژگی های غالب جمله را برجسته می کند.

3. **CLS Token:** در مدل BERT، توکن [CLS] (که در ابتدای جمله اضافه می شود) به عنوان نماینده جمله استفاده می شود. بردار خروجی مرتبط با این توکن، خلاصه ای از کل جمله است و می تواند برای بسیاری از تسک ها به کار گرفته شود.

روش های بهینه سازی Pooling با استفاده از Sentence-BERT:

Sentence-BERT یکی از مدل هایی است که برای استخراج نمای کلی جملات بهینه سازی شده است. این مدل از معماری BERT پایه استفاده کرده، اما تغییراتی در فرآیند تنظیم دقیق اعمال کرده است:

1. **هدف Sentence-BERT:** تبدیل هر جمله به یک بردار معنادار در فضای برداری، به گونه ای که جملات با معنای مشابه، به بردارهای نزدیک به هم تبدیل شوند.

2. کاربرد: **Pooling** در این مدل، از روش‌های مختلف Pooling مانند mean pooling یا max pooling استفاده می‌شود تا جمله به یک نمایه برداری ثابت تبدیل شود.

3. مزایا: استفاده از این تکنیک باعث می‌شود مدل بتواند به طور مستقیم در تسک‌هایی مانند مقایسه معنایی جملات، طبقه‌بندی متون و بازیابی اطلاعات استفاده شود.