



## پاسخ مسئله ۱.

(الف)

- با توجه به اطلاعات داده شده، در ابتدا تعداد positive pair ها را به دست می آوریم:

$$\# \text{ of positive pairs} = 2 + 3 + 6 \times 4 + 2 + 3 = 34$$

دقت کنید که با قرارگیری هر کدام از کلمات اول و دهم به عنوان کلمه مرکزی ۲ عدد positive pair، و با قرارگیری هر کدام از کلمات دوم و نهم به عنوان کلمه مرکزی ۳ عدد positive pair ایجاد خواهد شد. برای هر کدام از کلمات دیگر نیز ۴ عدد positive pair خواهیم داشت. حال دقت کنید که برای هر positive pair، ۵ عدد negative sample داریم و در نتیجه تعداد کل training sample ها از رابطه زیر به دست می آید:

$$\# \text{ of training samples} = 6 \times \# \text{ of positive pairs} = 204$$

- اگر target word در جایگاه  $i$  ام corpus قرار گرفته باشد (کلمه  $i$  ام corpus باشد) و طول corpus را با  $T$  نمایش دهیم، آنگاه تعداد input word های مورد استفاده برای پیش بینی آن از رابطه زیر محاسبه می شود:

$$\# \text{ of input words} = \min(i - 1, 4) + \min(T - i, 4)$$

- طبق اطلاعات مسئله،  $d_{model}$  برابر با ابعاد embedding ورودی و  $h$  برابر تعداد attention head ها است. پس ابعاد هر head برابر است با:

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

در نتیجه ماتریس های projection یعنی  $W_i^Q$ ،  $W_i^K$ ، و  $W_i^V$  برای  $i \in \{1, \dots, h\}$  دارای ابعاد زیر خواهند بود:

$$W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d_{model} \times d_k}$$

این تنظیمات باعث می شوند که عملیات محاسباتی روی بخش های کوچک تر از فضای برداری اصلی انجام شوند که باعث می شود محاسبات ماتریسی سریع تر انجام شوند، قابلیت پردازش موازی بهبود پیدا کند و کاهش حافظه مورد نیاز برای ذخیره ماتریس های projection اتفاق بیفتد.

• **صحیح.** در معماری اصلی Transformer، از fixed sinusoidal positional encodings استفاده می‌شود که به صورت از پیش تعریف شده و بدون نیاز به یادگیری اعمال می‌شوند. این کدگذاری‌های موقعیتی به منظور اضافه کردن اطلاعات ترتیبی به ورودی‌ها به کار می‌روند، زیرا مکانیزم Self-Attention ذاتاً نسبت به ترتیب توکن‌ها Permutation-Equivariant است.

در مقابل، مدل BERT از learnable positional embeddings استفاده می‌کند. این embeddings به صورت پارامترهای قابل یادگیری تعریف می‌شوند و طی فرایند training بهینه‌سازی می‌شوند. در نتیجه، برخلاف Transformer که از کدگذاری‌های سینوسی ثابت بهره می‌گیرد، BERT اطلاعات موقعیتی را به صورت پارامترهای مستقل و قابل یادگیری در طول آموزش به دست می‌آورد.

• **غلط.** افزایش تعداد attention head ها در مدل Transformer همیشه باعث بهبود عملکرد نمی‌شود. این افزایش تعداد می‌تواند به استخراج ویژگی‌های متفاوت کمک کند، اما معایب زیر را نیز ممکن است ایجاد کند:

- هزینه محاسباتی بالاتر: افزایش تعداد head ها، نیاز به محاسبات بیشتری دارد.
- خطر بیش‌برازش (Overfitting): مدل ممکن است به داده‌های آموزشی وابستگی زیادی پیدا کند و عملکرد آن روی داده‌های جدید کاهش یابد.
- نیاز به حافظه بیشتر: تعداد بیشتر head ها به معنای ذخیره ماتریس‌های بیشتر است که نیاز به حافظه بیشتری دارد.

## پاسخ مسئله‌ی ۲.

الف) به سادگی می‌توان نوشت:

$$\begin{aligned}
 -\frac{\partial \log P(w_o|w_t)}{\partial \mathbf{v}_{w_o}} &= -\frac{\partial \log P(w_o|w_t)}{\partial P(w_o|w_t)} \cdot \frac{\partial P(w_o|w_t)}{\partial \mathbf{v}_{w_o}} \\
 &= -\frac{1}{P(w_o|w_t)} \cdot \frac{\partial \left( \frac{\exp(\mathbf{u}_{w_t}^T \mathbf{v}_{w_o})}{\sum_{k \in \mathbf{v}} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_k)} \right)}{\partial \mathbf{v}_{w_o}} \\
 &= -\frac{1}{P(w_o|w_t)} \cdot \frac{\mathbf{u}_{w_t} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_{w_o}) \sum_{k \in \mathbf{v}} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_k) - \mathbf{u}_{w_t} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_{w_o})^2}{\left( \sum_{k \in \mathbf{v}} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_k) \right)^2} \\
 &= -\frac{1}{P(w_o|w_t)} \cdot \mathbf{u}_{w_t} \left( \frac{\exp(\mathbf{u}_{w_t}^T \mathbf{v}_{w_o})}{\left( \sum_{k \in \mathbf{v}} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_k) \right)} - \left( \frac{\exp(\mathbf{u}_{w_t}^T \mathbf{v}_{w_o})}{\sum_{k \in \mathbf{v}} \exp(\mathbf{u}_{w_t}^T \mathbf{v}_k)} \right)^2 \right) \\
 &= -\frac{1}{P(w_o|w_t)} \cdot \mathbf{u}_{w_t} (P(w_o|w_t) - P(w_o|w_t)^2) \\
 &= \mathbf{u}_{w_t} (P(w_o|w_t) - 1)
 \end{aligned}$$

ب) در جدول زیر اثر مقادیر مختلف  $c$  مقایسه شده‌اند.

اندازه پنجره ( $c$ )	تاثیر روی وابستگی‌های معنایی	پیچیدگی محاسباتی
۱	مدل تنها به همسایه‌های بسیار نزدیک توجه می‌کند و روی روابط محلی تمرکز دارد. این تنظیم می‌تواند برای داده‌هایی با ساختار محلی قوی مناسب باشد. با این حال، روابط طولانی‌مدت نادیده گرفته می‌شوند.	محاسبات سریع‌تر و بهینه‌تر انجام می‌شود زیرا تعداد جفت‌های مثبت و منفی کمتر است و حافظه کمتری نیاز دارد.
۵	یک تعادل بین روابط محلی و جهانی ایجاد می‌کند. مدل قادر به یادگیری ارتباطات معنی‌دار بین کلمات نزدیک و دورتر است. این اندازه معمولاً برای داده‌های عمومی مناسب است.	هزینه محاسباتی متوسط است. حافظه و زمان پردازش بیشتر از $c = 1$ اما کمتر از $c = 100$ است.
۱۰۰	تمرکز بر روابط جهانی و طولانی‌مدت دارد. این اندازه برای داده‌هایی که روابط دور برد دارند (مانند متون طولانی یا مقالات علمی) مناسب است، اما نویز بیشتری به مدل اضافه می‌کند.	نیاز به حافظه و زمان پردازش بالاتر دارد. تعداد نمونه‌های آموزشی به‌طور چشمگیری افزایش می‌یابد که منجر به بار محاسباتی زیاد می‌شود.

جدول ۱: مقایسه تاثیر اندازه پنجره در مدل skip-gram

### پاسخ مسئله‌ی ۳.

۱. : فرض کنید دنباله ورودی شامل بردارهای کلمه به شکل زیر است:

$$X = \{x_1, x_2, \dots, x_N\}$$

که  $N$  طول دنباله و  $D$  ابعاد هر بردار کلمه است. مکانیزم خودتوجه وزن توجه را برای هر کلمه در دنباله بر اساس روابط بین تمام کلمات در دنباله محاسبه می‌کند. عملیات خودتوجه به صورت زیر فرمول‌بندی می‌شود:

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

برای نمایش این عملیات به عنوان یک شبکه کاملاً متصل در قالب یک ماتریس، می‌توان بردارهای کلمه  $x_i$  را در کنار هم قرار داد تا یک ماتریس  $X'$  با ابعاد  $N \times D$  تشکیل شود که هر ردیف آن یک بردار کلمه را نشان می‌دهد. در این صورت خواهیم داشت:

$$Q = X'W^Q, \quad K = X'W^K, \quad V = X'W^V$$

حال، ماتریس توجه  $A$  را می‌توان به صورت زیر (با استفاده از ماتریس  $X'$  و ترانپوز آن) محاسبه کرد:

$$A = X'(X')^T$$

این ضرب ماتریسی، روابط بین تمام بردارهای کلمه در دنباله ورودی را نشان می‌دهد. درنهایت، برای نگاشت توالی کامل ورودی به یک بردار خروجی با همان ابعاد، می‌توان ضرب ماتریس دیگری را بین  $A$  و یک ماتریس وزن قابل یادگیری  $W$  به دست آورد:

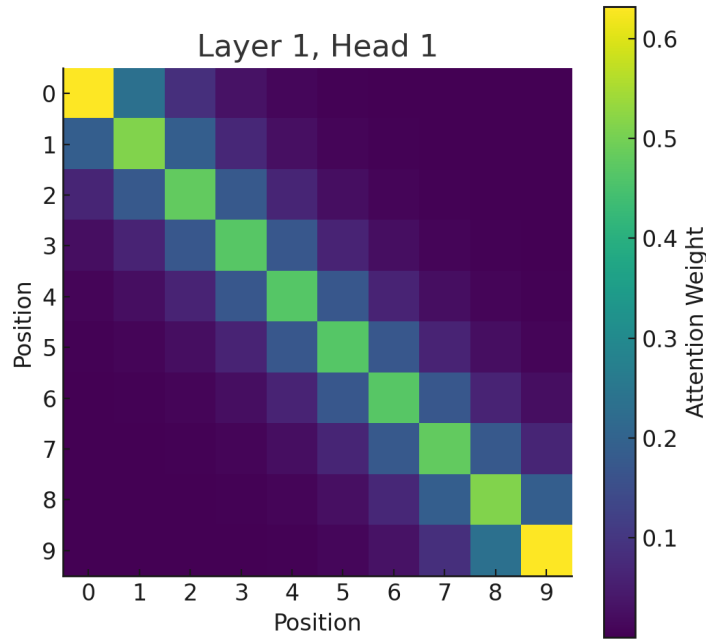
$$\text{Output} = \text{Softmax}(AW)$$

عبارت بالا خروجی یک شبکه تمام‌متصل با ماتریس وزن  $W$  و تابع فعال‌سازی Softmax است.

۲. دقت کنید که در شبکه به دست‌آمده در قسمت قبل، ماتریس  $W$  به شکل  $N \times N$  است (با فرض حفظ ابعاد یکسان) و شامل  $O(N^2)$  پارامتر می‌باشد. این نمایش شبکه کاملاً متصل، دارای  $O(N^2 D^2)$  پارامتر خواهد بود، زیرا  $N^2$  درایه در ماتریس توجه  $A$  وجود دارند و هر کدام آن‌ها حاصل ضرب بردارهای  $D$  بعدی است.

۳. در مکانیزم خودتوجه، هر کلمه در دنباله ورودی به همه کلمات دیگر توجه می‌کند، اما تمام توجه‌ها به یک اندازه مهم نیستند. برای دستیابی به این هدف، وزن توجه با استفاده از یک تابع Softmax محاسبه می‌شود که منجر به توزیع توجه پراکنده می‌شود، به طوری که تنها چند کلمه توجه قابل توجهی دریافت می‌کنند.

هر موقعیت در دنباله به همه موقعیت‌ها توجه می‌کند، اما قدرت توجه توسط ضرب نقطه‌ای تعیین می‌شود که منجر به یک ماتریس پراکنده می‌شود. در این ماتریس، بیشتر وزن‌های توجه نزدیک به صفر هستند. این پراکندگی از طریق تابع softmax به دست می‌آید که بر مرتبط‌ترین کلمات تأکید کرده و کلمات کمتر مرتبط را سرکوب می‌کند.



شکل ۱: ماتریس وزن‌ها در مکانیزم خودتوجه

- بلوک‌های سایه‌دار وزن‌های توجه غیرصفر را نشان می‌دهند که بیانگر اهمیت توجه است.
- بلوک‌های مورب نشان‌دهنده توجه به یک کلمه (توجه به خود) هستند که ممکن است به دلیل اهمیت توجه به یک کلمه صفر نباشد.
- بلوک‌های خارج از مورب نشان‌دهنده توجه به کلمات دیگر در دنباله هستند، جایی که بیشتر وزن‌های توجه به دلیل پراکندگی نزدیک به صفر هستند.

این ساختار پراکنده با به اشتراک‌گذاری پارامترها به مکانیزم خودتوجه اجازه می‌دهد تا وابستگی‌های بین کلمات را در توالی ورودی به‌طور مؤثر ثبت کند و درعین حال قابلیت پردازش محاسباتی را حفظ نماید.

۴. مکانیزم Self-Attention به دلیل ساختار خود، نسبت به ترتیب ورودی‌ها Permutation-Equivariant است؛ به این معنا که اگر ترتیب سطرهای ورودی جابه‌جا شود، خروجی نهایی نیز مطابق همان جابه‌جایی تغییر خواهد کرد. این ویژگی زمانی صادق است که مکانیزم Positional Encoding استفاده نشده باشد. در اینجا این رفتار را به طور دقیق‌تر بررسی می‌کنیم.

فرض کنید ورودی به صورت یک ماتریس  $X \in \mathbb{R}^{N \times D}$  باشد، که در آن  $N$  تعداد توکن‌ها و  $D$  نیز Embedding Dimension است. اکنون فرض کنید یک جایگشت سطرها با تابع

$$\pi : \{1, 2, \dots, N\} \rightarrow \{1, 2, \dots, N\}$$

اعمال شده است و ترتیب سطرهای  $X$  را به  $X' \in \mathbb{R}^{N \times D}$  تبدیل کرده است؛ به طوری که  $X'_{\pi(i)} = X_i$  برقرار می‌باشد. برای ماتریس ورودی  $X'$ ، بردارهای Query، Key، و Value به صورت زیر محاسبه می‌شوند:

$$Q' = X'W^Q, \quad K' = X'W^K, \quad V' = X'W^V.$$

به دلیل خطی بودن ضرب ماتریسی و اعمال جایگشت، رابطه زیر برقرار است:

$$Q'_{\pi(i)} = Q_i, \quad K'_{\pi(i)} = K_i, \quad V'_{\pi(i)} = V_i$$

اکنون ضرب درونی ماتریس‌های  $Q'$  و  $K'$  به صورت  $Q'K'^T$  خواهد بود. در اینجا، جایگشت در سطرها و ستون‌ها به طور همزمان اعمال می‌شود:

$$(Q'K'^T)_{\pi(i),\pi(j)} = (QK^T)_{i,j}$$

با اعمال Softmax روی هر سطر، ضرایب توجه نیز به همان ترتیب سطرها تحت جایگشت قرار می‌گیرند:

$$\text{Softmax}\left(\frac{Q'K'^T}{\sqrt{d_k}}\right)_{\pi(i),\pi(j)} = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)_{i,j}$$

بنابراین وزن‌های به‌دست‌آمده در Self-Attention مستقیماً با جایگشت ورودی مطابقت دارند. محاسبه خروجی نهایی با استفاده از بردارهای وزن‌دار شده  $V'$  به صورت زیر انجام می‌شود:

$$Y' = \text{Attention}(Q', K', V') = \text{Softmax}\left(\frac{Q'K'^T}{\sqrt{d_k}}\right)V' \implies Y'_{\pi(i)} = Y_i$$

و این همان چیزی است که می‌خواستیم نشان دهیم.

## پاسخ مسئله‌ی ۴.

۱. • ورودی Encoder: [۳۲, ۲۰۴۸, ۱۰۲۴] (سایز بچ، طول دنباله، ابعاد Embedding)
- خروجی Self-Attention در Encoder (هر سر): [۳۲, ۲۰۴۸, ۱۹۲]
- خروجی Self-Attention ادغام شده در Encoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی لایه FeedForward اول در Encoder: [۳۲, ۲۰۴۸, ۵۱۲]
- خروجی لایه FeedForward دوم در Encoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی Encoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- ورودی Decoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی Self-Attention در Decoder (هر سر): [۳۲, ۲۰۴۸, ۱۹۲]
- خروجی Self-Attention ادغام شده در Decoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی Cross-Attention در Decoder (هر سر): [۳۲, ۲۰۴۸, ۱۹۲]
- خروجی Cross-Attention ادغام شده در Decoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی لایه FeedForward اول در Decoder: [۳۲, ۲۰۴۸, ۵۱۲]
- خروجی لایه FeedForward دوم در Decoder: [۳۲, ۲۰۴۸, ۱۰۲۴]
- خروجی نهایی مدل: [۳۲, ۲۰۴۸, ۳۰۰۰۰]

۲. فرضیات داده شده:

- طول دنباله ورودی: ۲۰۴۸
- سایز واژگان ورودی: ۳۰۰۰۰
- بعد بُردار کلمه (Embedding Vector): ۷۶۸
- تعداد بلوک‌های Encoder: ۸
- تعداد بلوک‌های Decoder: ۱۲
- تعداد سرهای توجه (Attention Heads): ۴
- سایز بچ: ۳۲
- تعداد لایه‌های FeedForward: ۲

### محاسبات Encoder

(آ) پارامترهای Self-Attention:

$$d_k, d_v : \frac{768}{4} = 192$$

$$W_Q, W_K, W_V : 3 \times 1024 \times 192 = 589,824$$

$$W_O : 768 \times 1024 = 786,432$$

(ب) کل پارامترهای Multi-Head Attention:

$$4(589,824) + 786,432 = 3,145,728$$

(ج) پارامترهای FeedForward:

$$W_1 : 1024 \times 512 + 512 = 524,800$$

$$W_2 : 512 \times 1024 + 1024 = 525,312$$

(د) مجموع پارامترهای هر بلوک:

$$3,145,728 + 524,800 + 525,312 = 4,195,840$$

(ه) پارامترهای تمام بلوکهای Encoder:

$$8 \times 4,195,840 = 33,566,720$$

### محاسبات Decoder

(آ) مشابه Encoder، Decoder دارای پارامترهای زیر است: (دقت کنید که Decoder یک لایه Cross-Attention نیز دارد که تعداد پارامترهای آن کاملاً مشابه لایه Self-Attention می‌باشد)

$$12 \times (4,195,840 + 3,145,728) = 88,098,816$$

مجموع پارامترهای اولیه و نهایی

(آ) پارامترهای Embedding: (برای لایه linear آخر که خروجی Decoder را به خروجی نهایی مدل تبدیل می‌کند)

$$1024 \times 30000 + 30000 = 30,750,000$$

(ب) مجموع کل پارامترها:

$$33,566,720 + 88,098,816 + 30,750,000 = 152,415,536$$

۳. فرض کنید ورودی جمله زیر باشد:

”گربه روی تشک خوابیده است.”

و هدف ترجمه به انگلیسی باشد:

”The cat is sleeping on the mat.”



مدل با مراحل زیر عمل می‌کند:

### (آ) جاسازی توکن‌ها (Token Embedding):

- هر کلمه به برداری با ابعاد  $1024$  نگاشت می‌شود. این بردارها به صورت مقادیر عددی شناور نمایش داده می‌شوند و هر بردار اطلاعات معنایی و دستوری کلمه را رمزگذاری می‌کند.
- به عنوان مثال، جمله ورودی "گره روی تشک خوابیده است." به تعدادی توکن تقسیم شده و سپس هر توکن به یک بردار  $1024$  بُعدی تبدیل می‌شود.
- شکل خروجی اولیه این مرحله به صورت  $[32, 2048, 1024]$  است که  $32$  نشان‌دهنده سائز batch،  $2048$  طول دنباله (sequence length) و  $1024$  ابعاد هر بردار است.
- در این مرحله، Positional Encoding نیز به بردارها اضافه می‌شود تا مدل بتواند ترتیب کلمات را در نظر بگیرد.

### (ب) پردازش Encoder:

- ابتدا، برای هر توکن ورودی، ماتریس‌های  $Q, K, V$  (پرسش، کلید و مقدار) محاسبه می‌شوند. این ماتریس‌ها ابعاد  $[32, 2048, 192]$  برای هر Head دارند.
- محاسبات Attention با استفاده از فرمول زیر انجام می‌شود:

$$\text{Attention}(q, k, v) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

- که در آن  $d_k = 192$  است. این عملیات برای یافتن میزان توجه به هر توکن دیگر در دنباله استفاده می‌شود.
- سپس، نتایج Head ها ادغام شده و پس از ضرب شدن در ماتریس  $W^O \in \mathbb{R}^{64 \times 1024}$  به ابعاد اولیه بازگردانده می‌شود ( $[32, 2048, 1024]$ ).
- نتایج از دو لایه FeedForward عبور می‌کنند:  $\sigma$  تابع فعال‌سازی است

$$\text{FFN}(X) = \sigma(\sigma(XW_1 + b_1)W_2 + b_2)$$

- لایه‌های FeedForward به ترتیب ابعاد  $[32, 2048, 512]$  و  $[32, 2048, 1024]$  را تولید می‌کنند که خروجی دومین لایه FeedForward همان خروجی نهایی Encoder است.

### (ج) پردازش Decoder:

- ابتدا Self-Attention مشابه Encoder برای توکن‌های هدف محاسبه می‌شود.
- سپس، Cross-Attention برای ارتباط بین خروجی Encoder و ورودی Decoder انجام می‌شود.
- فرمول Cross-Attention به صورت زیر است:

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

که در آن  $Q$  از Decoder و  $K, V$  از Encoder محاسبه می‌شوند.

- مشابه Encoder، نتایج از دو لایه FeedForward عبور کرده و بردارهای نهانی نهایی تولید می‌شوند.

(د) پیش‌بینی کلمات:

- خروجی Decoder به یک ماتریس خطی با ابعاد  $[1024 \times 30000]$  ضرب می‌شود که در آن ۳۰۰۰۰ تعداد کل واژگان است.
- با استفاده از تابع Softmax، احتمالات توزیع روی واژگان محاسبه می‌شود:

$$p(w) = \frac{e^{z_w}}{\sum_{v \in V} e^{z_v}}$$

که در آن  $z_w$  امتیاز (logit) مربوط به واژه  $w$  است.

- واژه‌ها به صورت Autoregressive تولید می‌شوند. این به این معنی است که هر کلمه به ترتیب پیش‌بینی شده و به ورودی اضافه می‌شود تا کلمه بعدی پیش‌بینی گردد.
- فرآیند پیش‌بینی تا رسیدن به توکن پایان جمله ([EOS]) یا تکمیل طول حداکثر ادامه می‌یابد.

## پاسخ مسئله‌ی ۵.

الف) مدل BERT در زمان آموزش اولیه خود از هدف (NSP) Next Sentence Prediction برای یادگیری روابط معنایی بین جملات متوالی استفاده می‌کند. در این فرآیند، مدل با جفت جملات آموزش داده می‌شود تا تشخیص دهد که آیا جمله دوم به‌طور منطقی به جمله اول متصل است یا خیر. این هدف به مدل کمک می‌کند تا قابلیت درک روابط متنی و ترتیب جملات را بیاموزد. با این حال، در برخی تسک‌های پردازش زبان طبیعی، نیاز به چنین اطلاعاتی وجود ندارد.

در مدل‌هایی مانند RoBERTa، هدف NSP حذف شده و تمرکز کامل بر روی Masked Language Modeling (MLM) قرار گرفته است. دلیل اصلی این تغییر آن است که بسیاری از وظایف پایین‌دستی، مانند طبقه‌بندی متن یا تشخیص احساسات، به درک معنایی کلی جملات نیاز دارند و نیازی به بررسی روابط ترتیبی بین جملات ندارند. حذف NSP باعث می‌شود منابع محاسباتی بیشتری به یادگیری ارتباطات معنایی میان کلمات اختصاص یابد و عملکرد مدل در این نوع وظایف بهبود پیدا کند.

همچنین، حذف NSP پیچیدگی آموزش را کاهش می‌دهد و باعث می‌شود مدل به‌طور کارآمدتری برای تسک‌های متنی مانند پاسخ به سوالات یا ترجمه آموزش داده شود. مدل‌هایی مانند RoBERTa با تکیه بر این استراتژی موفق به دستیابی به نتایج بهتری در مقایسه با BERT در بسیاری از معیارهای استاندارد ارزیابی عملکرد شده‌اند.

ب) برای بهبود عملکرد BERT در وظایف پایین‌دستی، روش‌هایی مانند Sentence-BERT (SBERT) معرفی شده‌اند که هدف آن‌ها افزایش توانایی مدل در درک معنایی جملات کامل است. برخلاف BERT که بر روی توکن‌های منفرد تمرکز دارد، SBERT از جملات به‌عنوان واحدهای معنایی مستقل استفاده می‌کند و بردارهای متراکم و معنایی برای جملات استخراج می‌کند.

در SBERT، به‌جای استفاده از خروجی‌های مربوط به توکن‌های منفرد، کل جمله به یک بردار تبدیل می‌شود که خلاصه‌ای از محتوای آن است. این روش معمولاً از تکنیک pooling برای تولید بردار نهایی جمله استفاده می‌کند. در این فرآیند، ابتدا بردارهای خروجی مربوط به تمام توکن‌های جمله از آخرین لایه مدل استخراج می‌شوند. سپس از یکی از روش‌های pooling مانند mean pooling، max pooling یا انتخاب بردار مربوط به توکن [CLS] برای ترکیب اطلاعات تمامی توکن‌ها به یک بردار استفاده می‌شود.

روش mean pooling با محاسبه میانگین بردارهای تمامی توکن‌ها، یک نمای کلی و متوازن از اطلاعات جمله ایجاد می‌کند. در مقابل، max pooling با انتخاب بیشترین مقدار در هر بُعد بردارها، بر ویژگی‌های قوی‌تر تاکید دارد. استفاده از توکن [CLS] نیز برای نمایندگی جمله محبوب است، زیرا این توکن به‌طور خاص در طول آموزش برای خلاصه‌سازی اطلاعات کل جمله آموزش دیده است.

برداری نهایی به‌دست‌آمده از pooling می‌تواند به‌عنوان ورودی برای مدل‌های پایین‌دستی مانند تطابق معنایی جملات یا خوشه‌بندی متن به کار رود. این رویکرد به‌ویژه برای تسک‌هایی که نیاز به مقایسه جملات دارند، مانند ارزیابی شباهت معنایی یا بازیابی اطلاعات، بسیار کارآمد است.

یکی دیگر از مزایای SBERT این است که به دلیل استفاده از بردارهای معنایی جمله، در مقایسه با BERT عملکرد بهتری در وظایف ارزیابی شباهت جملات دارد. همچنین، این روش برای تسک‌هایی مانند پاسخ به سوالات و طبقه‌بندی متن بسیار موثرتر است، زیرا نیاز به پردازش جداگانه توکن‌ها را از بین می‌برد.

به‌طور کلی، تنظیم دقیق BERT با استفاده از روش‌هایی مانند SBERT، امکان یادگیری بردارهای معنایی دقیق‌تر را فراهم می‌کند و باعث می‌شود مدل در تسک‌های پایین‌دستی عملکرد بهتری داشته باشد.