



نام و نام خانوادگی:

رضا قربانی پاچی

شماره دانشجویی:

403206565

تمرین پنجم درس یادگیری ماشین

# سوال 1

(الف)

## الف) قسمت 1

وقتی سائز پنجره برابر 2 است هر کلمه می تواند با دو کلمه قبل و بعدش جفت بشه. مثلاً برای جفت های اول و آخری این تعداد به 2 و 3 جفت میرسه در نهایت خواهم داشت

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

برای هر جفت مثبت  $k = 5$  تا صفت می خواهم داشت

$$\Rightarrow \text{مجموعه جفت های مثبت} = k \times \text{تعداد جفت های مثبت} = 5 \times 34 = 170$$

$$\Rightarrow \text{مجموعه جفت های منفی} = 204 + 170 = 374$$

## الف) قسمت 2

در مدل CBOW اگر اندازه پنجره 4 باشد (یعنی context window size برابر 4)، تعداد کلمات ورودی (input words) برای پیش بینی هر target word با  $\text{ext} = 8$  خواهد بود شامل 4 کلمه از سمت چپ و 4 کلمه از سمت راست هدف است

## الف) قسمت 3

اگر ابعاد embedding ورودی  $d_{\text{model}}$  باشد و تعداد  $h$  head وجود داشته باشد، هر head بهی برای  $\frac{d_{\text{model}}}{h}$  خواهد داشت. بنابراین ابعاد ماتریس های  $w_v$ ،  $w_k$ ،  $w_q$  و  $w_r$  می خواهد بود head برابر با

$$w_q: d_{\text{model}} \times \frac{d_{\text{model}}}{h}$$

$$w_k: \sim \times \sim$$

$$w_v: \sim \times \sim$$

تنظیم این پارامترها به کاهش پیچیدگی محاسباتی کمک می کند زیرا تعداد عملیات کاهش یافته و محاسبات به صورت مکانیکی روی head ها توزیع می شود

(ب)

ب- ۱)

صمیم است. در مدل positional encoding، BERT به صورت trainable embedding است و در طول فرآیند آموزش بهینه می‌شود.  
در مقابل، در معماری اصلی positional encoding، Transformer به صورت ثابت و سینوسی طراحی شده است.

---

ب- ۲)

خلاف است. افزایش تعداد attention heads همیشه باعث بهبود عملکرد نمی‌شود. اگر تعداد head ها بیش از حد افزایش یابد، ممکن است منجر به افزایش پیچیدگی محاسباتی و مشکلات overfitting شوند. همچنین، تاثیر افزایش تعداد layers بر کیفیت مدل و تنظیمات و فرآیند دیگر مانند اندازه embedding دارد.

## سوال 2

(الف)

الف) سوال (الف)  $P(\omega_o | \omega_t)$  به صورت زیر مشتق می‌گردد

$$P(\omega_o | \omega_t) = \frac{e^{\mu_{\omega_t}^T v_{\omega_o}}}{\sum_{k \in V} e^{\mu_{\omega_t}^T v_k}} \xrightarrow{\text{ارگونم}} \log P(\omega_o | \omega_t) = \mu_{\omega_t}^T v_{\omega_o} - \log \sum_{k \in V} e^{\mu_{\omega_t}^T v_k}$$

$$\begin{aligned} - \frac{\partial \log P(\omega_o | \omega_t)}{\partial v_{\omega_o}} &= - \frac{\partial \mu_{\omega_t}^T v_{\omega_o}}{\partial v_{\omega_o}} - \frac{\partial \log e^{\mu_{\omega_t}^T v_k}}{\partial v_{\omega_o}} \\ &= - \mu_{\omega_t} - \frac{1}{\sum_{k \in V} e^{\mu_{\omega_t}^T v_k}} \cdot \frac{\partial}{\partial v_{\omega_o}} \sum e^{\mu_{\omega_t}^T v_k} \end{aligned}$$

برای  $k = \omega_o$  داریم

$$- \frac{\partial \log P(\omega_o | \omega_t)}{\partial v_{\omega_o}} = - \mu_{\omega_t} + \mu_{\omega_t} \cdot P(\omega_o | \omega_t)$$

(ب)

: c=1

- یک پنجره کوچک روابط بافت محلی را ثبت می‌کند. به این معنا که:
  - مدل بر روی جفت کلماتی که مستقیماً در کنار یکدیگر ظاهر می‌شوند، تمرکز می‌کند.
  - نمایش‌هایی که یاد گرفته می‌شوند، بیشتر روابط نحوی (مانند ساختار دستوری یا وابستگی کلمات) را منعکس می‌کنند (مثلاً "گربه" و "است" یا "یک").
- انتظار: تعبیه‌های کلمات (word embeddings) به خوبی وابستگی‌های کوتاه‌برد را ثبت می‌کنند، اما ممکن است در مدل‌سازی روابط معنایی گسترده‌تر ناتوان باشند.

:C=5

- یک پنجره متوسط هم بافت محلی و هم برخی روابط معنایی را ثبت می‌کند.

○ کلماتی که در یک ساختار شبیه به جمله ظاهر می‌شوند، مد نظر قرار می‌گیرند و تعادلی بین الگوهای نحوی و معنایی ایجاد می‌شود.

- انتظار: تعبیه‌های کلمات، اطلاعات نحوی و معنایی را به طور همزمان نمایش می‌دهند و آن‌ها را برای بسیاری از وظایف پردازش زبان طبیعی کاربردی می‌سازد.

C=100:

- یک پنجره بزرگ روابط معنایی گسترده را ثبت می‌کند و بر روی هم‌وقوعی کلمات در دامنه‌های بزرگ‌تر تمرکز دارد.

○ کلماتی که از نظر موضوعی به هم مرتبط هستند، اما لزوماً از نظر نحوی نزدیک نیستند، بیشتر در نظر گرفته می‌شوند (مثلاً "گره" و "حیوان").

- انتظار: تعبیه‌ها وابستگی‌های معنایی بلندبرد را ثبت می‌کنند که برای وظایفی که به درک معنایی نیاز دارند (مانند مدل‌سازی موضوعی یا اندازه‌گیری شباهت معنایی) مفید است. با این حال، ممکن است جزئیات ظریف نحوی از دست برود.

## قسمت 1

در مکانیزم Self-Attention، دنباله‌ی ورودی که شامل جاسازی کلمات (یا توکن‌ها) است از طریق سه ماتریس متفاوت پردازش می‌شود:

1. پرسش‌ها: (Q) این ماتریس نمایانگر کلمه یا توکنی است که در حال حاضر مورد توجه قرار گرفته است. برای هر کلمه در دنباله، یک بردار پرسش محاسبه می‌شود تا اهمیت آن نسبت به سایر توکن‌ها مشخص شود.

2. کلیدها: (K) این بردارها ویژگی‌هایی از هر توکن در دنباله را نشان می‌دهند که ممکن است با یک پرسش مطابقت داشته باشد. کلیدها برای محاسبه امتیازات توجه با پرسش‌ها جفت می‌شوند.

3. مقادیر: (V) این ماتریس حاوی اطلاعات واقعی توکن‌هاست که بر اساس میزان اهمیت (امتیازات توجه) وزن‌دهی می‌شوند.

مراحل در فرم ماتریسی:

1. محاسبه‌ی امتیازات توجه: (Dot Product Attention)

امتیازات توجه از طریق ضرب نقطه‌ای بین ماتریس پرسش (Q) و ترانپوز ماتریس کلید ( $K^T$ ) محاسبه می‌شود. به صورت ریاضی:

$$Scores = Q \cdot K^T$$

این عملیات تعیین می‌کند که هر پرسش چقدر به هر کلید در دنباله مرتبط است. اگر حاصل ضرب نقطه‌ای بالا باشد، نشان‌دهنده‌ی ارتباط قوی است.

2. مقیاس‌دهی با  $\sqrt{D}$ :

بعد از بردار کلید (D) برای مقیاس‌دهی امتیازات توجه استفاده می‌شود. این مقیاس‌دهی با تقسیم امتیازات بر  $\sqrt{D}$  از بزرگ شدن بیش از حد مقادیر در هنگام کار با ابعاد بالا جلوگیری می‌کند:

$$Scaled\ Scores = \frac{Q \cdot K^T}{\sqrt{D}}$$

### 3. تابع Softmax:

امتیازات مقیاس شده از طریق تابع Softmax نرمال سازی می شوند تا به احتمال تبدیل شوند. این مرحله تضمین می کند که جمع وزن های توجه برای هر بردار پرسش برابر ۱ باشد:

$$Attention\ Weights = Softmax\left(\frac{Q \cdot K^T}{\sqrt{D}}\right)$$

این وزن ها نشان دهنده ی میزان اهمیت نسبی هر کلمه در دنباله برای یک پرسش خاص هستند.

### 4. وزن دهی مقادیر:

وزن های توجه سپس در ماتریس مقدار (V) ضرب می شوند. این مرحله اطلاعات همه ی توکن ها را بر اساس ارتباطشان با پرسش فعلی ترکیب می کند:

$$Y = Softmax\left(\frac{Q \cdot K^T}{\sqrt{D}}\right) \cdot V$$

این فرمول بندی ماتریسی به صورت کارآمد توجه را برای تمام توکن ها در دنباله به طور همزمان محاسبه می کند. نتیجه (Y) دنباله ای از بردارهاست، که هر بردار خروجی ترکیب وزن دهی شده ای از بردارهای مقدار ورودی است. وزن ها بر اساس شباهت بین پرسش و کلیدها و نرمال سازی توسط Softmax تعیین می شوند.

این فرآیند هسته ی اصلی مکانیزم Self-Attention است که به مدل هایی مانند Transformer امکان می دهد تا وابستگی های بلندمدت و روابط متنی در دنباله ها را به خوبی درک کنند.

## قسمت 2)

تحلیل ابعاد به صورت گام به گام:

1. ماتریس‌های ورودی:

○ ماتریس‌های  $Q$  (پرسش‌ها) و  $K$  (کلیدها) ابعادی برابر با  $N \times D$  دارند:

▪  $N$ : طول دنباله که تعداد توکن‌های موجود در دنباله ورودی است.

▪  $D$ : ابعاد بردارهای کلمه (یا اندازه جاسازی هر توکن).

2. ماتریس کلید ترانهاده‌شده:

○ ماتریس ترانهاده‌شده  $K^T$  ابعادی برابر با  $D \times N$  دارد. ترانهاده کردن، سطرها و ستون‌های  $K$  را جابجا می‌کند.

3. ضرب ماتریس:  $(Q \cdot K^T)$

○ ضرب  $Q(N \times D)$  در  $K^T(D \times N)$  ماتریسی با ابعاد  $N \times N$  ایجاد می‌کند.

○ هر عنصر در این ماتریس  $N \times N$  نشان‌دهنده امتیاز توجه بین دو توکن در دنباله است (مثلاً اینکه توکن  $i$  چقدر به توکن  $j$  توجه می‌کند).

4. نرمال‌سازی: Softmax

○ تابع Softmax به صورت عنصری روی هر سطر ماتریس  $N \times N$  اعمال می‌شود. این کار اطمینان می‌دهد که مجموع وزن‌های توجه برای هر توکن برابر ۱ باشد. اندازه ماتریس همچنان  $N \times N$  باقی می‌ماند.



5. ضرب در مقادیر ( $V$ ):

- ماتریس امتیازات توجه ( $N \times N$ ) سپس در ماتریس  $V$  که ابعاد  $N \times D$  دارد ضرب می‌شود. این مرحله کمک می‌کند تا سهم وزنی تمام توکن‌ها برای هر پرسش ترکیب شود.
- خروجی این عملیات ماتریسی با ابعاد  $N \times D$  خواهد بود که با ابعاد ماتریس اصلی پرسش یا مقدار یکسان است.

تحلیل پارامترها:

تعداد کل پارامترهای مورد استفاده در محاسبات Self-Attention به صورت زیر تقسیم می‌شود:

1. ماتریس‌های  $Q$ ،  $K$  و  $V$  هر کدام دارای  $N \times D$  عنصر هستند، زیرا از ماتریس‌های وزنی یاد گرفته شده اعمال شده روی دنباله ورودی تولید می‌شوند.

- در مجموع، این سه ماتریس  $3 \times N \times D$  پارامتر ایجاد می‌کنند.

2. هنگام محاسبه  $Q \cdot K^T$ ،  $N^2$  عملیات ضرب نقطه‌ای انجام می‌شود که هر کدام نیاز به  $D$  ضرب و جمع دارند.

- این مرحله  $N^2 \times D$  عملیات به خود اختصاص می‌دهد.

3. ضرب ماتریس امتیازات توجه ( $N \times N$ ) در  $V$  ( $N \times D$ ) نیاز به  $N^2 \times D$  عملیات اضافی دارد.

توضیحات بالا پیچیدگی محاسبات را برای یک گام (Attention Head) مشخص می‌کند. حال اگر بخواهیم تمام پارامترهای مدل را در نظر بگیریم:

1. هر (Attention Head) شامل سه ماتریس یادگیری‌پذیر است  $W_Q$ ،  $W_K$ ، و  $W_V$ .

- هر کدام از این ماتریس‌ها ابعاد  $D \times D$  دارند.

- تعداد کل پارامترها در این سه ماتریس  $3 \times D^2$  است.

2. حال اگر  $N \times D$  توکن را پردازش کنیم و بخواهیم محاسبات مربوط به هر توکن را برای همه توکن‌های دنباله انجام دهیم:

○  $D^2$  پارامترهای ماتریس‌ها ضرب در  $N^2$  محاسبات برای توجه زوجی بین تمام توکن‌ها می‌شود.

در نتیجه:

$$\text{Total Complexity} = O(N^2 D^2)$$

## قسمت 3

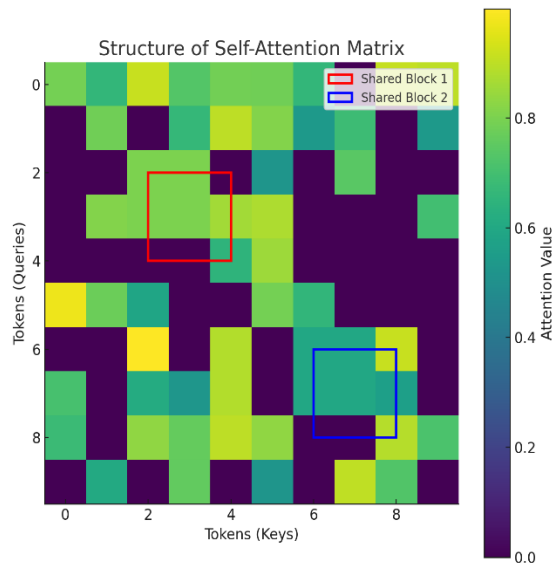
### 1. اشتراک پارامترها در ماتریس Self-Attention

- در مکانیزم Self-Attention، بلوک‌های ماتریس تعامل ( $Q \cdot K^T$ ) اغلب پارامترهای مشترک دارند:
  - پارامترهای مشترک به دلیل وجود وزن‌های یادگیری‌پذیر ( $W_Q, W_K, W_V$ ) برای همه توکن‌ها هستند.
  - به عنوان مثال، اگر دو توکن وزن‌های توجه مشابهی برای سایر توکن‌ها داشته باشند، بلوک‌های مربوط به تعامل آن‌ها در ماتریس توجه یکسان خواهند بود.

### 2. ساختار ماتریس Sparse

- در یک ماتریس Sparse:
  - بسیاری از عناصر به دلیل عملیات Softmax و محدودیت مقادیر کوچک به صفر نزدیک یا دقیقاً صفر هستند.
  - این ساختار خلوت باعث کاهش حجم محاسبات می‌شود، زیرا عملیات تنها بر روی عناصر غیر صفر انجام می‌شود.
- در ساختار ماتریس:
  - بلوک‌هایی که تعامل قوی دارند (وزن توجه بالا)، مقادیر غیر صفر خواهند داشت.
  - بلوک‌های دیگر که تعامل کمی دارند یا اهمیتی ندارند، به صفر تبدیل می‌شوند.

### 3. نمودار ساختار ماتریس



این نمودار ساختار یک ماتریس Self-Attention را نمایش می‌دهد.

- هر عنصر در ماتریس نشان‌دهنده تعامل (وزن توجه) بین یک جفت توکن است.
- بلوک‌های قرمز و آبی نشان‌دهنده بخش‌هایی از ماتریس هستند که پارامترهای مشترک دارند.
- عناصر نزدیک به صفر (نمایش داده شده با رنگ‌های تیره) ناشی از تابع Softmax هستند، که خلوت بودن (Sparsity) ماتریس را نشان می‌دهند.

## قسمت 4)

### اثر حذف Positional Encoding

#### 1. بدون: Positional Encoding

- اگر Positional Encoding حذف شود، مدل قادر به تشخیص موقعیت کلمات در دنباله ورودی نخواهد بود.
- در نتیجه، خروجی لایه Attention تنها بر اساس روابط محتوایی (Semantic) بین کلمات محاسبه می‌شود، بدون در نظر گرفتن ترتیب آن‌ها.
- به عنوان مثال:
  - دنباله ورودی "The cat sat on the mat" :
  - ترتیب تغییر یافته "Mat the on sat cat the" :
  - خروجی برای هر دو دنباله مشابه خواهد بود، زیرا مدل موقعیت کلمات را در نظر نمی‌گیرد.

#### 2. ترتیب مجدد دنباله‌ها:

- اگر ترتیب دنباله‌ها تغییر کند، خروجی به ترتیب دنباله جدید خواهد بود. این به این معناست که مدل همچنان خروجی معادل تولید می‌کند، اما این خروجی دیگر ترتیب معنایی اصلی را حفظ نمی‌کند.

اثر حذف بردارهای ورودی

1. اگر بردارهای ورودی حذف شوند:

- بدون ورودی‌های معنایی (Semantic Embeddings)، مدل هیچ اطلاعاتی درباره کلمات یا توکن‌های ورودی ندارد.
- در این حالت، خروجی Attention به طور کامل بی‌معنی خواهد بود، زیرا هیچ محتوای معنایی برای محاسبه ارتباط وجود ندارد.
- تمام توجه‌ها و خروجی‌ها ممکن است به صورت تصادفی یا صفر تولید شوند.

2. وابستگی به: Positional Encoding

- حتی اگر Positional Encoding به تنهایی باقی بماند، بدون اطلاعات معنایی بردارهای ورودی، مدل نمی‌تواند تصمیم‌گیری مناسبی درباره ترتیب و معنای کلمات داشته باشد.

## سوال 4

### Encoder:

- Input Shape: (BatchSize,SequenceLength,HiddenDimension)=(32,2048,768)
- Key, Query, Value Sizes:  
Each size = (HiddenDimension,HiddenDimension)=(768,768)
- Self-Attention Layer in Encoder Block:
  - Input Shape: (32,2048,768)
  - Output Shape: (32,2048,768)
- Feed-Forward Network in Encoder Block:
  - Input Shape: (32,2048,768)
  - Intermediate Output: (32,2048,384)
  - Final Output: (32,2048,768)
- Output Dimensions per Encoder Block:  
(32,2048,768)
- Total Number of Encoder Blocks: 8
  - Final Output Shape after Encoder:  
(32,2048,768)

### Decoder:

- Input Shape: (BatchSize,TargetSequenceLength,HiddenDimension)=(32,1024,768)
- Key, Query, Value Sizes:  
Each size = (HiddenDimension,HiddenDimension)=(768,768)
- Self-Attention Layer in Decoder Block:
  - Input Shape: (32,1024,768)
  - Output Shape: (32,1024,768)
- Encoder-Decoder Attention in Decoder Block:
  - Input Shape: (32,1024,768)

- Output Shape: (32,1024,768)
- Feed-Forward Network in Decoder Block:
  - Input Shape: (32,1024,768)
  - Intermediate Output: (32,1024,384)
  - Final Output: (32,1024,768)
- Output Dimensions per Decoder Block: (32,1024,768)
- Total Number of Decoder Blocks: 12
  - Final Output Shape after Decoder: (32,1024,768)

(ب)

Parameter Calculation for the Architecture:

#### 1. Embedding Layers (Encoder + Decoder):

- Embedding Matrix Dimensions:

$$\text{VocabSize} \times \text{HiddenSize} = 30,000 \times 768 = 23,040,000$$

For both encoder and decoder:

$$23,040,000 \times 2 = 46,080,000 \text{ parameters}$$

#### 2. Encoder Blocks:

Each encoder block has two main components:

- Self-Attention Layer:
  - Query, Key, Value Matrices:  $\text{HiddenSize} \times \text{HiddenSize} = 768 \times 768 = 589,824$

Total for all three matrices:

$$589,824 \times 3 = 1,769,472 \text{ parameters}$$

Output projection matrix:

$$\text{HiddenSize} \times \text{HiddenSize} = 768 \times 768 = 589,824 \text{ parameters}$$



Total Self-Attention Parameters:

$$1,769,472 + 589,824 = 2,359,296 \text{ parameters}$$

- Feed-Forward Network (FFN):

- First layer:

$$\text{HiddenSize} \times \text{FFN\_Size} = 768 \times 384 = 294,912$$

- Second layer:  $\text{FFN\_Size} \times \text{HiddenSize} = 384 \times 768 = 294,912$

- Total FFN Parameters:

$$294,912 + 294,912 = 589,824$$

- Total Parameters per Encoder Block:

$$2,359,296 + 589,824 = 2,949,120 \text{ parameters}$$

- Total for 8 Encoder Blocks:

$$2,949,120 \times 8 = 23,592,960 \text{ parameters}$$

### 3. Decoder Blocks:

Each decoder block has three main components:

- Self-Attention Layer:

Same as Encoder: 2,359,296 parameters

- Encoder-Decoder Attention Layer:

Same as Self-Attention: 2,359,296 parameters

- Feed-Forward Network (FFN):

Same as Encoder: 589,824 parameters

- Total Parameters per Decoder Block:

$$2,359,296 + 2,359,296 + 589,824 = 5,308,416 \text{ parameters}$$

- Total for 12 Decoder Blocks:

$$5,308,416 \times 12 = 63,700,992 \text{ parameters}$$

#### 4. Output Layer:

- Fully connected layer mapping hidden size to vocabulary size:  $\text{HiddenSize} \times \text{VocabSize} = 768 \times 30,000 = 23,040,000$  parameters

#### Total Parameters:

- Embedding Layers: 46,080,000
- Encoder Blocks: 23,592,960
- Decoder Blocks: 63,700,992
- Output Layer: 23,040,000

#### Grand Total:

$46,080,000 + 23,592,960 + 63,700,992 + 23,040,000 = 156,413,952$  parameters

(ج)

تسک: ترجمه جمله "I love programming." از انگلیسی به فارسی

1. ورودی به انکودر:

- جمله ورودی "I love programming.":
- تبدیل جمله به توکن‌ها (اندیس‌ها) [350, 200, 101]: نمایانگر "I", "love", "programming"
- شکل ورودی به لایه جاسازی (Embedding Layer):  
 $(\text{BatchSize}, \text{SequenceLength}) = (3, 1)$
- بعد از جاسازی: هر توکن به برداری با بُعد 768 تبدیل می‌شود:  
 $(\text{BatchSize}, \text{SequenceLength}, \text{HiddenDimension}) = (768, 3, 1)$

## 2. کدگذاری موقعیت: (Positional Encoding)

- به هر توکن ورودی، اطلاعات مربوط به موقعیت آن در جمله اضافه می‌شود (مثلاً "I" اولین کلمه است)

3. عبور از بلاک‌های انکودر (8 بلاک):

در هر بلاک انکودر:

### 1. لایه توجه چندسری: (Multi-Head Self-Attention)

○ ورودی (768,3,1) :

○ سه ماتریس Query (Q)، Key (K) و Value (V) ایجاد می‌شود :

▪ شکل هر ماتریس (768,3,1) :

○ مکانیزم توجه محاسبه می‌شود Attention(Q,K,V)

○ خروجی: همان شکل ورودی (768,3,1) (1,3,768)

### 2. عملیات نرمال سازی: (Add & Norm)

○ خروجی توجه به ورودی اولیه اضافه می‌شود و نرمال سازی انجام می‌گیرد.

○ نتیجه (768,3,1) :

### 3. شبکه پیش‌خور: (Feed-Forward Network)

○ لایه اول (افزایش بُعد):

ورودی : (768,3,1) , خروجی : (384,3,1)

○ لایه دوم (کاهش بُعد):

ورودی : (384,3,1) , خروجی : (768,3,1)

### 4. عملیات نرمال سازی دوم: (Add & Norm)

○ خروجی شبکه پیش‌خور به نتیجه توجه اضافه می‌شود و نرمال سازی انجام می‌گیرد.

○ خروجی (768,3,1) :

4. خروجی نهایی انکودر:

- بعد از عبور از 8 بلاک انکودر، شکل نهایی خروجی :  
 $(BatchSize, SequenceLength, HiddenDimension) = (768, 3, 1)$
- این خروجی شامل اطلاعات غنی شده درباره روابط و معانی کلمات در جمله انگلیسی است که به دیکودر ارسال می شود.

1. ورودی به دیکودر:

- دیکودر ترجمه را به صورت مرحله به مرحله تولید می کند.
- ورودی اولیه به دیکودر، تنها توکن شروع (START) است.  
مثال:

Input Tokens: [1]

- این توکن از طریق لایه جاسازی (Embedding Layer) به یک بردار 768 بعدی تبدیل می شود :  
شکل ورودی به دیکودر :  $(BatchSize, TargetSequenceLength, HiddenDimension) = (768, 1, 1)$

2. عبور از بلاک های دیکودر (12 بلاک):

در هر بلاک دیکودر:

1. لایه توجه چندسری (Self-Attention):

- این لایه تنها روی دنباله خروجی فعلی دیکودر (ترجمه تولیدشده تاکنون) تمرکز می کند.
- سه ماتریس Query (Q)، Key (K) و Value (V) ایجاد می شود :
  - شکل هر ماتریس  $(768, 1, 1)$  برای مرحله اول.
- مکانیزم توجه محاسبه می شود  $Attention(Q, K, V)$
- (Masked Attention): برای اطمینان از اینکه دیکودر تنها به خروجی های قبلی خود دسترسی دارد، از ماسک استفاده می شود.
- خروجی: همان شکل ورودی  $(768, 1, 1)$

2. عملیات نرمال سازی اول: (Add & Norm)

- خروجی توجه به ورودی اولیه اضافه می شود و نرمال سازی انجام می گیرد (768,1,1)

3. توجه انکودر-دیکودر: (Encoder-Decoder Attention)

- در این مرحله، دیکودر روی خروجی انکودر تمرکز می کند.
- ماتریس های Query (Q)، Key (K) و Value (V) به ترتیب از دیکودر و انکودر ایجاد می شوند :

▪ Q از دیکودر (768,1,1) :

▪ K و V از انکودر (768,3,1) :

- مکانیزم توجه محاسبه می شود  $\text{Attention}(Q, K, V)$  :

- خروجی (768,1,1)

4. عملیات نرمال سازی دوم: (Add & Norm)

- خروجی توجه انکودر-دیکودر به خروجی قبلی اضافه شده و نرمال سازی می شود (768,1,1) :

5. شبکه پیش خور: (Feed-Forward Network)

- لایه اول (افزایش بُعد) (768,1,1) → (384,1,1) :

- لایه دوم (کاهش بُعد) (384,1,1) → (768,1,1) :

6. عملیات نرمال سازی سوم: (Add & Norm)

- خروجی شبکه پیش خور به نتیجه قبلی اضافه شده و نرمال سازی می شود (768,1,1) :

3. پیش بینی کلمه خروجی:

- خروجی دیکودر به یک لایه کاملاً متصل (Fully Connected Layer) داده می شود که اندازه آن برابر با تعداد کلمات واژگان (Vocabulary Size) است:

•  $\text{VocabularySize} \rightarrow \text{HiddenDimension} = 30,000 \rightarrow 768$

- خروجی این لایه احتمال هر کلمه در واژگان را نشان می دهد. کلمه با بیشترین احتمال انتخاب می شود :

Predicted Token: [7] ("من" در زبان فارسی)

4. تکرار مراحل برای کلمه بعدی:

- توکن پیش‌بینی‌شده (مثلاً [7]) به دنباله ورودی دیکودر اضافه می‌شود.
- دیکودر مجدداً مراحل بالا را تکرار می‌کند تا کلمه بعدی پیش‌بینی شود :

○ گام 1 " → [1] :من"

○ گام 2 " → [7,1] :عاشق"

○ گام 3 " → [50,7,1] :برنامه‌نویسی"

○ گام 4 <END> → [200,50,7,1] :

5. خروجی نهایی دیکودر:

- خروجی نهایی دیکودر، ترجمه کامل جمله است:  
"من عاشق برنامه‌نویسی هستم."

## سوال 5

(الف)

حذف مرحله «پیش‌بینی جمله بعدی (NSP)» از پیش‌پردازش BERT در مطالعات مختلف (مانند معماری RoBERTa) نشان داده است که می‌تواند در بسیاری از وظایف پایین‌دستی (Downstream Tasks) عملکرد مدل را بهبود دهد. دلایل اصلی چنین بهبودی را می‌توان به صورت زیر خلاصه کرد:

1. سادگی و محدودیت: NSP  
در BERT اصلی، NSP با یک کار ساده (تشخیص «آیا جمله دوم از نظر توالی طبیعی به جمله اول متصل است یا نه») آموزش می‌بیند. این کار، مدل را وادار می‌کند تا نوعی تشخیص «تصادفی بودن یا نبودن» بین دو جمله انجام دهد که الزاماً در بسیاری از وظایف زبانی پیچیده مفید نیست. بنابراین بخشی از ظرفیت مدل صرف یادگیری الگویی می‌شود که در بسیاری از کاربردهای عملی نقش کلیدی ندارد.
2. کاهش نویز و تمرکز بیشتر بر بافت کلمه‌ها:  
در رویکرد NSP، جفت جمله‌های ساختگی (Random) ممکن است باعث ایجاد سیگنال‌های آموزشی پرت (Noise) شود و مدل برای تشخیص همین سیگنال از داده‌های ماسک‌شده (MLM) منحرف شود. حذف NSP به مدل اجازه می‌دهد توجه و ظرفیتش را به شکل متمرکزتری روی پیش‌بینی توکن‌های ماسک‌شده بر اساس بافت (Context) صرف کند و در نتیجه نمایش (Representation) کلی بهتری از زبان به دست آورد.
3. الگوی نامناسب اتصال جمله‌ها:  
روش NSP معمولاً با انتخاب جمله‌هایی از همان سند یا از اسناد متفاوت صورت می‌گیرد. اما در داده‌های واقعی، نحوه گذار از یک جمله به جمله بعدی می‌تواند بسیار متنوع‌تر از «متصل/نامتصل» ساده باشد. به عبارت دیگر، بسیاری از متون دنیای واقعی (مثل پاراگراف‌های طولانی، مباحث چندبخشی و غیره) به روش NSP پوشش مناسبی داده نمی‌شوند. حذف NSP باعث می‌شود مدل درگیر چنین الگوی محدودی نشود.
4. بهبود عملکرد در وظایفی که نیازمند انسجام درون جمله‌ای هستند:  
بسیاری از وظایف پایین‌دستی (مانند تحلیل احساس، طبقه‌بندی جمله، استخراج موجودیت‌ها و غیره) بیش از آنکه به رابطه میان دو جمله احتیاج داشته باشند، نیازمند فهم عمیق‌تر محتوای تک جمله یا مجموع جملات هستند. حذف NSP به مدل امکان می‌دهد تا از لایه‌های بیشتری برای فهم بهتر کلمات و جملات استفاده کند و در نتیجه، در این وظایف ساده‌تر به دقت بالاتری برسد.
5. یافته‌های تجربی:  
در معماری RoBERTa و برخی مدل‌های دیگر، با حذف NSP مشاهده شد که مدل در اکثر وظایف GLUE و سایر

بنچمارک‌ها نتایج بهتری کسب می‌کند. این یافته‌های تجربی نشان می‌دهد که هدف NSP آن‌چنان که تصور می‌شد برای فراگیری روابط میان جمله‌ای مؤثر نبوده و حتی می‌تواند موجب هدررفت ظرفیت مدل شود.

در مجموع، اگرچه ایده NSP در ابتدا برای بهبود درک روابط سطح بالا بین جملات در نظر گرفته شده بود، اما در عمل، شکل ساده و دوحالتی آن (پیوسته بودن یا نبودن جملات) کمک چندانی به نمایش بهتر زبان نمی‌کند و حتی ممکن است مانع شود که مدل به صورت عمیق‌تری بر خود کلمات و بافت درونی جملات تمرکز کند. از این رو، حذف NSP می‌تواند در کاربردهای واقعی، به ویژه وظایف پایین دستی که بیشتر بر درک محتوای تک جمله یا چند جمله پیوسته متمرکزند، عملکرد را بهبود بخشد.

## (ب)

برای رسیدن به «درک کلی از جمله» در هنگام تنظیم دقیق (Fine-tuning) مدل‌هایی مانند BERT یا RoBERTa و به کارگیری آن‌ها در تسک‌های پایین دستی (مثلاً طبقه‌بندی جملات، تحلیل احساس، بازیابی متون و غیره)، معمولاً نیاز داریم که ورودی (یک جمله یا عبارت) را به یک بردار نهایی تبدیل کنیم؛ برداری که بتواند به عنوان نماینده یا چکیده معنایی کل جمله ورودی عمل کند. به این فرایند اغلب «نمایش جمله» (Sentence Representation) گفته می‌شود. در ادامه، چند روش متداول برای دستیابی به چنین نمایشی توضیح داده شده است:

### 1. استفاده از بردار [CLS] روش پیش فرض BERT

در معماری BERT، توکن ویژه‌ای به نام [CLS] در ابتدای دنباله قرار داده می‌شود. در خروجی آخرین لایه ترانسفورمر (یا یکی از لایه‌های آخر بسته به طراحی)، بردار متناظر با [CLS] به نوعی نقش «خلاصه‌کننده» یا «نماینده کل دنباله» را دارد. در بسیاری از سناریوهای طبقه‌بندی، از همین بردار به عنوان ورودی یک شبکه کوچک (مثلاً یک لایه Dense) استفاده می‌شود تا پیش‌بینی نهایی انجام گیرد.

- **مزیت:** پیاده‌سازی ساده و سازگاری با طراحی اولیه BERT.
- **چالش:** لزوماً تضمین نمی‌شود که بردار [CLS] همیشه «بهترین» یا «تنها» راه برای نمایش کل جمله باشد. در عمل، گاهی مشاهده می‌شود که روش‌های دیگر (مثلاً تجمیع همه توکن‌ها) نتایج بهتری در برخی تسک‌ها ارائه می‌کنند.



## 2. استفاده از ( Pooling ) تجميع بردارهای همه توکن‌ها)

در این روش، به جای اتکا به بردار خاصی مانند [CLS]، تمام توکن‌های خروجی از آخرین لایه یا چند لایه آخر مدل را در نظر می‌گیریم. سپس با نوعی عملگر pooling (مانند میانگین، حداکثر یا ترکیبی از آن‌ها) یک بردار واحد را به دست می‌آوریم. به این ترتیب، نمود هر توکن (و به تبع آن بخش‌های مختلف جمله) در نمایش نهایی لحاظ می‌شود.

### • Average Pooling:

$$Representation = \frac{1}{N} \sum_{i=1}^n \frac{1}{n} h_i$$

### • Max Pooling:

$$Representation[d] = \max_{1 \leq i \leq n} (h_i[d])$$

که  $h_i[d]$  مؤلفه dام بردار نهفتگی توکن i است و در پایان، ماکس در هر بعد گرفته می‌شود.

### • مزیت:

- توجه به تمام توکن‌ها و در نتیجه دریافت اطلاعات جامع‌تر از جمله.
- سادگی پیاده‌سازی (فقط یک لایه ساده‌ی Pooling و قابل استفاده روی هر لایه از خروجی (برای مثال آخرین لایه یا میانگین چند لایه)).

### • چالش:

- برخی مواقع ممکن است نویزهای موجود در بخشی از جمله بر تجميع کلی غلبه کنند (خصوصاً در روش Max Pooling).
- انتخاب مناسب نوع Pooling و لایه/لایه‌هایی که از آن Pool می‌گیریم، نیازمند آزمون و خطا است.
- 

## 3. مدل‌های تخصصی برای نمایش جمله (مانند Sentence-BERT)

**Sentence-BERT** یا به اختصار **SBERT**، معماری‌ای مبتنی بر BERT است که مشخصاً برای نمایش سطح جمله طراحی شده و اغلب به صورت Siamese Network (دوقلو) آموزش می‌بیند. در این روش:

1. دو جمله (یا بیشتر) به صورت موازی از مدل گذر داده می‌شوند.

2. در انتهای هر شاخه، یک لایه Pooling معمولاً Average Pooling یا CLS Pooling روی خروجی آخرین لایه اعمال می‌شود تا بردار نهایی جمله به دست آید.

3. برای آموزش، از زبان‌های مبتنی بر شباهت بردارها (مثلاً کازین شباهت یا کنتراستی) استفاده می‌شود تا مدل بیاموزد جمله‌های «مشابه» بردارهای نهایی نزدیکی داشته باشند و جمله‌های «نامشابه» بردارهای دورتری داشته باشند.

- **مزیت:**

- آموزش صریحاً بر اساس نمایش سطح جمله انجام می‌شود، بنابراین در وظایف متناظر (مثلاً بازیابی مبتنی بر شباهت جمله، خوشه‌بندی جملات، یا پرسش-پاسخ‌های کوتاه) کارایی بالاتری دارد.
- سرعت بالاتر در انجام مقایسه جملات (چون هر جمله یک بار به یک بردار تبدیل می‌شود و نیازی به pairwise comparison در خود مدل نیست)

- **چالش:**

- راه‌اندازی و آموزش اولیه آن نیازمند داده‌های مناسب (مثلاً جفت‌جملات مشابه/نامشابه) و کمی پیچیدگی بیشتر نسبت به استفاده مستقیم از BERT دارد.

#### 4. لایه‌های اضافه و معماری‌های ترکیبی

در برخی موارد، برای به‌دست‌آوردن نمایش بهتر سطح جمله، یک لایه یا بلاک اضافی به خروجی مدل اضافه می‌شود که کارش «تجمیع و تلفیق بردارهای همه توکن‌ها» است. این بلاک می‌تواند شامل مکانیسم توجه (Attention) یا مکانیزم‌های پیشرفته‌تر باشد تا کاستی‌های pooling ساده برطرف شود. به عنوان مثال، ممکن است:

1. ابتدا با یک یا چند لایه ترانسفورمر اضافی، بردارهای توکن‌ها دوباره با هم تعامل کنند.
2. در انتها از یک لایه pooling (مثلاً Mean Pooling) یا حتی مکانیزمی شبیه attention-pooling استفاده شود. این رویکرد گاهی در معماری‌های پیشرفته‌تر دیده می‌شود تا از تمام اطلاعات در سطوح مختلف مدل بیشترین بهره گرفته شود.

## جمع‌بندی

- درک کلی از جمله را می‌توان با استفاده از بردار خاص [CLS] (روش پیش‌فرض)، یا انواع Pooling (میانگین، حداکثر و ...)، یا مدل‌های تخصصی نظیر Sentence-BERT به دست آورد.
  - انتخاب بهترین راهکار، وابسته به نوع تسک پایین‌دستی است. برای مثال، در تسک‌های «تشخیص شباهت جملات» یا «جستجوی معنایی» روش‌هایی مثل SBERT معمولاً نتایج بهتری می‌دهند. اما در تسک‌های ساده‌تر (مثل طبقه‌بندی دودویی جملات) گاهی همان بردار [CLS] یا یک Pooling ساده، به اندازه کافی خوب عمل می‌کند.
  - **تنظیم دقیق (Fine-tuning)** نیز نقش مهمی دارد. حتی یک روش ساده مانند بردار [CLS] هم در صورتی که به‌درستی روی داده‌های تسک هدف (Task-specific data) تنظیم (Fine-tune) شود، می‌تواند عملکرد قابل قبولی داشته باشد. اما گاه با کمی تلاش بیشتر در طراحی مرحله نمایندگی جمله (مثلاً با استفاده از Pooling یا روش‌های تناوبی)، می‌توان بهبود محسوسی در عملکرد نهایی به دست آورد.
- در نتیجه، برای آنکه مدلی مانند BERT یا RoBERTa بتواند «درک کلی از جمله» ارائه دهد، لازم است یا از مکانیزم‌های مناسب در لایه خروجی (نظیر Pooling) استفاده کنیم یا از مدل‌هایی بهره بگیریم که به‌طور اختصاصی برای تولید بردار جملات آموزش دیده‌اند (همچون Sentence-BERT).