

лекция 6. Transformers. GPT. BERT. T5.

Transformers

1 Self-Attention

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \in \mathbb{R}^{T \times H}, \quad T - \text{длина под-ти}$$

$$Q, K \in \mathbb{R}^{T \times \tilde{H}}, \quad V \in \mathbb{R}^{T \times H}$$

$$\text{Self-Attention} = \text{Attention}(E \cdot W_Q, E \cdot W_K, E \cdot W_V),$$

$$W_Q, W_K \in \mathbb{R}^{H \times \tilde{H}}, \quad W_V \in \mathbb{R}^{H \times H} - \text{обучаемые матрицы}$$

2 Multi head Self-Attention

$$\text{head}_i = \text{Attention}\left(\underbrace{Q W_{i,Q}}_{H \times \tilde{H}}, \underbrace{K W_{i,K}}_{H \times \tilde{H}}, \underbrace{V \cdot W_{i,V}}_{H \times H}\right),$$

$$\text{MultiHeadAtt}(Q, K, V) = \text{concat}\left[\underbrace{\text{head}_1, \dots, \text{head}_n}_{T \times n\tilde{H}}\right] \cdot W^O \in \mathbb{R}^{T \times H} \quad \begin{array}{l} W^O \in \mathbb{R}^{n\tilde{H} \times H} \\ \text{— изначальная} \\ \text{размерность} \end{array}$$

На практике:

1) Нам не хочется заводить 3 матрицы и умножать их отдельно

$$\text{на } E \in \mathbb{R}^{T \times H}$$

↓ внутренняя раз-ть h , $Q, K, V \in \mathbb{R}^{H \times h}$

$$W \in \mathbb{R}^{H \times 3h} \Rightarrow [QKV] = E \cdot W \in \mathbb{R}^{T \times 3h}$$

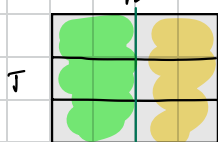
2) Теперь разделим их

$$Q, K, V = \text{split}([QKV])$$

$$\Rightarrow Q, K, V \in \mathbb{R}^{T \times h}$$

3) Теперь нужно разделить матрицы между головками, но опять будет не оптимально умножать.

$$\text{Подходит } A = Q \times K^T ?$$



Тогда очевидно не подходит.

reshape $Q, K, V: T \times h \rightarrow n_heads \times T \times h'$

$$A = \frac{Q K^T}{\sqrt{h'}} \in n_heads \times T \times T$$

$$A \leftarrow A + M$$

$$4) A = \text{Softmax}(A)$$

Также часто на этом этапе применяется Dropout

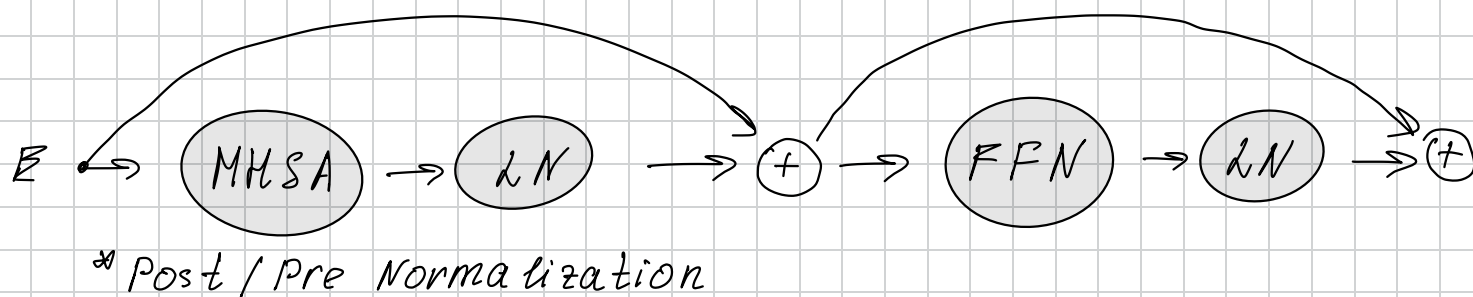
$$5) y = A \cdot V \in n_heads \times T \times h'$$

reshape $n_heads \times T \times h' \rightarrow T \times n_heads \times h' \rightarrow T \times h$

$$6) y = y \cdot W \rightarrow T \times M$$

Тут тоже обычно стоит dropout.

3 Transformer block



4 Position Encoding

$PE_{pos} \in \mathbb{R}^d$, pos-позиция в пос-ти

$$PE(pos, 2i) = \sin\left(\frac{pos}{f_i}\right) \quad f_i = 10000^{\frac{2i}{d}}$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{f_i}\right) \quad i = 1, \frac{d}{2}$$

Такие эмбединги обладают очень хорошими свойствами, что $pos+k$ эмбединг

является линейной функцией от pos эмбединга:

$$\begin{aligned} PE(pos+k, 2i) &= \sin\left(\frac{pos+k}{f_i}\right) = \sin\left(\frac{pos}{f_i}\right)\cos\left(\frac{k}{f_i}\right) + \cos\left(\frac{pos}{f_i}\right)\sin\left(\frac{k}{f_i}\right) = \\ &= \text{Linear}[PE(pos, 2i), PE(pos, 2i+1)] \end{aligned}$$

Transfer Learning

Рассказ про GPT, BERT, ERNIE, XLNET, Roberta

