

1. 词表大小、word $v \times h$ (17)

2. 设单句语长度为 L

1711 BERT 结构中 - Embedding

Token Embeddings 大小为 $\text{Token} \times h$ ①

segment embeddings 大小为 $\text{seg} \times h$ ②

Position Embeddings 大小为 $\text{pos} \times h$ ③

embeddings = ① + ② + ③ 为 embeddings $L \times h$.

1711 - 10 $P_n \in L \cdot W_{h \times 1} \quad e \in L \cdot b_{h \times 1}$ ④

Embedding 层中各词向量的求和

$$E_{P_n} = ① + ② + ③ + ④ = L \times h + 1 \times h + L \times h + h \times 1 + h \times 1 \\ = (2L + 3)h. \quad \dots \dots (17)$$

3 BERT-self-attention

$Q_w (h \times h) \quad Q_b (h \times 1)$ ⑤

$K_w (h \times h) \quad K_b (h \times 1)$ ⑥

$V_w (h \times h) \quad V_b (h \times 1)$ ⑦

Output P_n . $O_w (h \times h) \quad O_b (h \times 1)$ ⑧

self-attention P_n 训练参数为:

$$S_{\frac{L}{h}} = 4 \times h \times h + 4 \times h \times 1 = 4 \times h \times (1+h) \dots (2)$$

$$= 4h(1+h)$$

4 layer norm (X embedding + X attention)

$$E_{LW}(h \times 1), E_{Lb}(h \times 1)$$

$$E_{L\frac{L}{h}} = 2 \times h = 2h \dots (3)$$

$$5. \text{Output} = \text{liner}(\text{gelu}(\text{liner}(x)))$$

$$\text{liner}(x) = X_{L \times h} \cdot W_{h \times 4h} + b_{4h \times 1} = X_{L \times 4h}$$

$$O_{GLW}(h \times 4h) \quad O_{GLb}(4h \times 1) \dots (1)$$

$$\text{liner}(\text{gelu}(\text{liner}(x))) = X_{L \times 4h} \cdot W_{4h \times h} + b_{h \times 1}$$

$$O_{LW}(4h \times h) \quad O_{Lb}(h \times 1) \dots (2)$$

$$\text{Output}_{\frac{L}{h}} = (1) + (2) = 4h \times h \times 2 + 4h \times 1 + h \times 1$$

$$= h(8h + 5) \dots (4)$$

6 layer norm (X forward + X attention)

$$L_{\frac{L}{h}} = 2h$$

$$L_{\frac{L}{h}} = 2h \dots (5)$$

7. pooler. dense.

$$P_{\cdot W}(h \times h) \quad P_{\cdot b}(h \times 1)$$

$$P_{\frac{1}{2}} = h(h+1) \quad , \quad - \quad - \quad - (b)$$

BERT 结构 - 层网 收敛 数.

$$B_{\frac{1}{2}} = (1) + (4 + \dots + 17)$$

$$\begin{aligned}
 &= (2L+3)h + 4(1+h)h + 2h + (5+8h)h + 2h + (h+1)h + vh \\
 &= (2L+3+4+4h+2 + 5+8h+2+h+1+v)h \\
 &= (2L+13h+v+15)h
 \end{aligned}$$