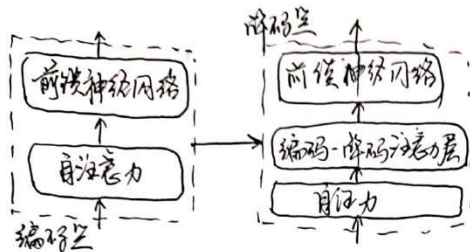


Transformer

① 67 编码器, 67 解码器

所有编码器结构相同, 但不共享参数



② 单词被编码成 512 维向量

1) 词嵌入只发生在最低层编码器的输入
高层编码器输入是低层编码器的输出
向量列表大小是超参, 一般是最大句子的长度

2) 位置编码

$$p_t = f(t) = \begin{cases} \sin(w_k \cdot t), & \text{if } i=2k \\ \cos(w_k \cdot t), & \text{if } i=2k+1 \end{cases}$$

$$w_k = \frac{1}{10000^{2k/d}}$$

频率沿向量维度减小, p_t 是一个包含每个频率的正弦余弦对

$$p_t = \begin{bmatrix} \sin(w_1 \cdot t) \\ \cos(w_1 \cdot t) \\ \sin(w_2 \cdot t) \\ \cos(w_2 \cdot t) \\ \vdots \\ \sin(w_d/2 \cdot t) \\ \cos(w_d/2 \cdot t) \end{bmatrix}$$

用论文加位置编码加到模型输入上

$$\phi'(w, t) = \phi(w, t) + p_t, \quad d_{w, t} = d_{p, t}$$

词嵌入更容易关注相对位置信息, 相对位置之间具有线性关系

$$e_t = E_{t, :} = \begin{bmatrix} \sin(\frac{t}{10000}) \\ \cos(\frac{t}{10000}) \\ \sin(\frac{t}{10000^2}) \\ \cos(\frac{t}{10000^2}) \\ \vdots \\ \sin(\frac{t}{10000^{d/2}}) \\ \cos(\frac{t}{10000^{d/2}}) \end{bmatrix}$$

$$f_m = \frac{1}{\lambda_m} = 10000 \frac{2m}{dm}$$

存在一线性投影 $T^k \in R^{d_m \times d_m}$, 对 $\forall t \in \{1, \dots, n-k\}$
矩阵位置偏移量 $k \in \{1, \dots, n\}$, 下列

$$T^{(k)} E_{t, :} = E_{t+k, :}$$

3) 自注意力机制

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V$$

X 在第 i 层是编码, 前后面是前一层的输出

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

1) 使用不同的 Q, K , 在不同空间投影, 增强表达能力, 提高泛化能力

2) 多头机制

多头扩展了模型关注不同位置的能力
给出了多个“表示子空间”

4) 残差网络

残差 + 归一化

每子层都有一残差连接

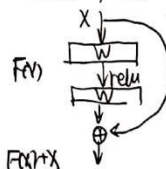
① 梯度消失/爆炸

② 网络退化

随着深度增加, 网络的表现是无增加至饱和, 然后迅速下降

对神经网络来说, 恒等映射并不容易拟合
(更深的网络不应该表现更差)

$$z^{(l)} = H(a^{(l)}) = a^{(l-1)} + F(a^{(l)})$$



1) 解决网络退化

2) 缓解梯度消失

3) 集成学习角度

4) 缓解梯度破碎

信息前后传播更顺畅