

BERT

① 特点: 1) 端到端, 中间层 24, 共 12 层
2) MLM

② mask 掉 15% 的词, 但用三种:

- 1) 80% 直接 mask
- 2) 10% 随机替换 因为全 mask 掉可能
- 3) 10% 保持不变 防止期间从未看到

③ NSP: 判断语句对是否连续

④ 模型架构

(1) 使用 Transformer 原始架构

(2) 输入表示:

$\begin{cases} \text{token embedding} \\ \text{segment embedding} \\ \text{position embedding: 可学习的} \end{cases}$

添加 [CLS], 用于分类任务

西关键创新

1) 应用 Cloze 任务, 去噪声编码

要求预测被 mask 的词

2) 缺点: 1. 没有见过被 mask 掉词

2. 每 batch 只预测了 15% token.
表明需要更多步骤才能收敛

3) NSP, 下一句预测

(4) 优点: ① 双向

② Cloze

③ NSP + MLM

④ 通用: 输入层输出层不用定制

⑤ 微调成本低

缺点: ① 随机 mask

② mask 实际预测没有过多影响

③ 收敛慢

④ 消耗巨大

5) 应用:

1) 命名实体识别任务本身

2) 句子级分类

3) 深层语义任务

4) 输入长度不是很长

Q: 为什么忽略点积

A: 维度 d 较高有较大方差, softmax 梯度较小, 缩放点积缓解这一问题

Q: 点积模型与加法模型

A: 复杂度差不多, 但点积可以用矩阵乘积效率更高

Q: multi Head 为什么有效

A: 更多特征表示方向, 类似 CNN 多通道选择, 无通过切头再分别进行 scale dot, 可以防止维度 d 不过大, 同时缩小 attention mask 影响

Q: FFN

A: FFN 中 σ ReLU 主要是非线性单元

Q: weight tying

A: 减少参数量

Q: GELU

A: 随机 erf , 输入 x 乘以 $-\frac{1}{2}$ mask, mask 依赖于输入
mask, m 服从正态分布

$$\text{GELU}(x) = x \cdot P(x \leq x) = x \cdot \Phi(x)$$

Q: GELU 优点

A: 缺乏随机因素 RELU 只用 0, 1,

Q: 优点 BERT

1) 双向 2) 并行化