

Finetune

1. 全量微调(Full Fine-tuning)

- 全量参数微调
- 模型结构和参数均进行调整
- 训练时间长，消耗资源大

2. PEFT(PEFT - Parameter-Efficient Fine-Tuning)

LoRA (Low-Rank Adaptation)

- 低秩适应
- 通过低秩分解更新权重
- 训练效率高

数学公式

$$W \in R^{r \times k} \text{ Lora } h = Wx + \triangle Wx = Wx + BAx \quad \dots$$

$$B \in R^{d \times r} A \in R^{r \times k} \quad r \leq \min(d, k)$$

- 低秩矩阵 $\triangle W$ 的计算
- 低秩矩阵 $\triangle W$ 的更新
- 低秩矩阵 $\triangle W$ 的正交化
- 低秩矩阵 $\triangle W$ 的计算公式
 - $\triangle W = \Sigma V^T r \Sigma$
 - $\Sigma = \sqrt{\lambda} I$

代码示例

- $A @ B @ B^T @ W @ 0$

Adapter

- Transformer模型的适配器
- Adapter
- Adapter

模型

Transformer模型的适配器

模型

Adapter模型

- 低秩矩阵 $d @ r$ (\$r \leq d\$)
- RELU或GELU
- $r @ d$

模型

模型 $Adapter(x) = x + W_{\text{up}} \cdot \text{ReLU}(W_{\text{down}} \cdot x)$

- $W_{\text{down}} \in R^{r \times d}$, $W_{\text{up}} \in R^{d \times r}$
- 低秩矩阵 $d @ r$

Bias-only

模型

- 偏置项(bias)的适配器

- 1%
1%
1%