

摘要

ELMo	深的双向 transformers
使用的是 RNN 的架构, 所以在应用	GPT 单向的左边的去预测未来
到下游任务的时候需要进行一些架构的调整	Bert 使用了双向的信息 同在下游任务只需要改上层 (类似 GPT)
	① 与哪些文章有关, 并且区别在哪 ② 好的点绝对精度 相对好的程度
	引言:
	feature-based ELMo (RNN 的架构)
	fine-tuning GPT 从左往右的架构 (transformer 架构)
	<ul style="list-style-type: none">• 带掩码的语言模型• 下一个句子的预测任务
贡献:	<ul style="list-style-type: none">• 双向信息的重要性• 微调调
	不是预测未来而是完形填空 cloze task

大量的没有标号的数据集上去训练比在只有少量有标号的数据上训练的效果更好

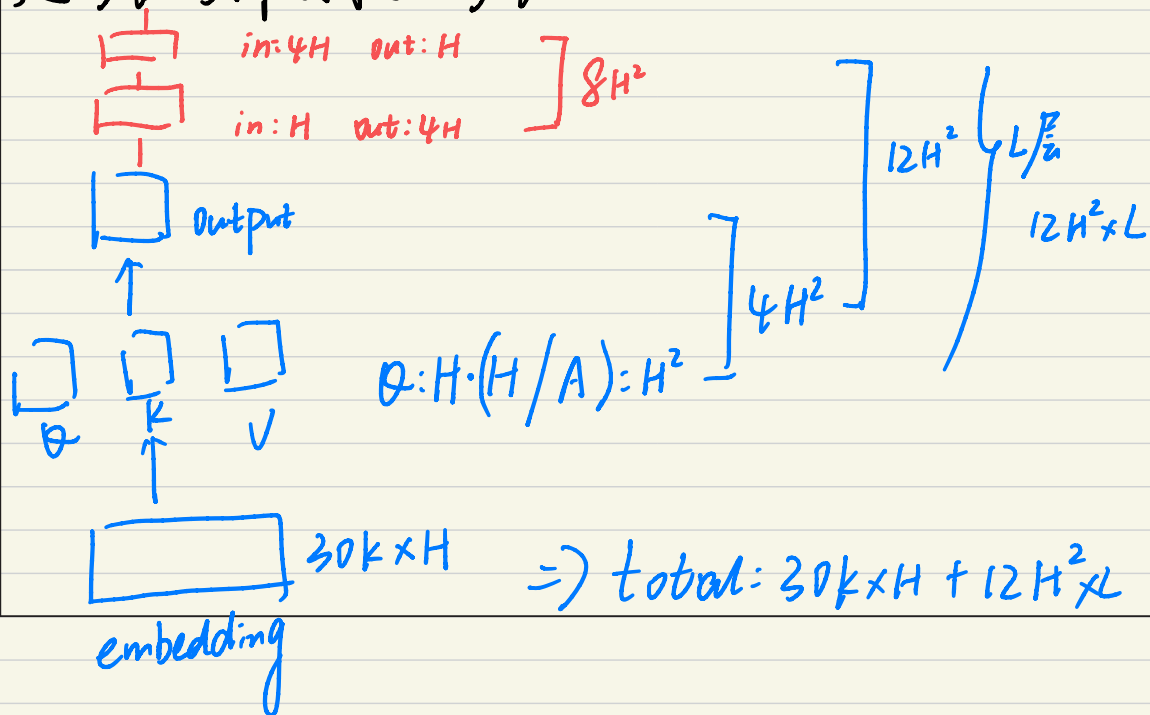
预训练: 在沒有标号的数据上进行训练

微调: 在下游任务使用有标号的数据进行微调

L, H, A Bert-base
 ↓ ↓
 hidden heads Bert-large

模型的宽度:

超参数转换成可学习参数的大小



word piece (看了序列) 30k

bert 输入是序列

transformer 输入是序列对

[CLS] [SEP]

bert 中 pos-emb 和 Seg-emb 都是通过训练学习得来的

MLM

Bert 更好为自然语言不方便生成