

模型设计报告

2024 春《机器学习》期末作业 —— 赛道二：日志异常检测

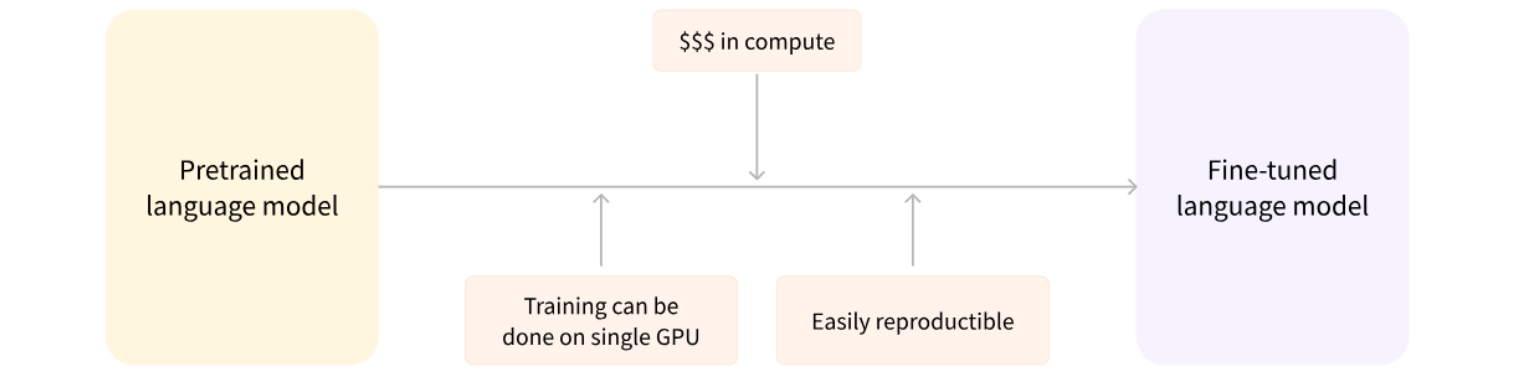
小组成员：

- 邓思阳 2021211352
- 梁骋 2021211221
- 朱宏明 2021211229

概述

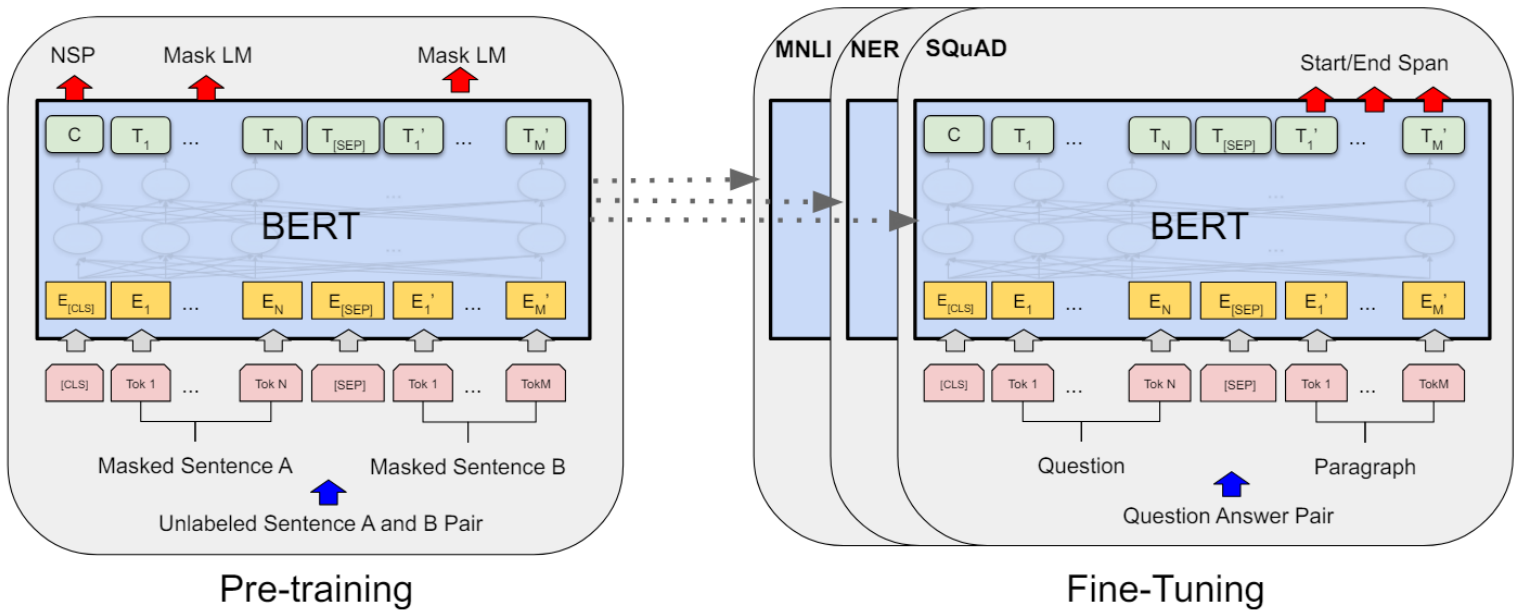
数据集为某系统运行过程中所产生的日志信息，为了实现系统自动报警能力，需要实现一个机器学习算法，利用日志信息判断系统是否有异常发生。我们收集到了多条日志信息，并对其进行了标注，标注分为两类：Normal 和 Anomalous。例如：`"log": "1135434878 2005.12.24 R67-M1-NC-C:J04-U01 RAS KERNEL r16=0xffffffff r17=0x4c00081f r18=0x44000000 r19=0x085f9780", "label": "Normal"`。

我们使用当下流行的大语言模型来讲解决日志是否异常这样一个文本二分类问题。基于 bert-base-cased 对其使用我们的数据集进行微调，得到可以完成指定日志是否异常的二分类任务。



模型结构

BERT 是一个 Transformer 模型，以自我监督的方式在大型英语数据语料库上进行预训练。这意味着它仅对原始文本进行预训练，没有人以任何方式标记它们（这就是为什么它可以使用大量公开可用的数据），并通过自动过程从这些文本生成输入和标签。



通过这种方式，模型可以学习英语语言的内部表示，然后可以使用该表示来提取对下游任务有用的特征：例如，如果有标记句子的数据集，则可以使用 BERT 生成的特征来训练标准分类器模型作为输入。

训练方法与测试结果

使用 Huggingface 提供的 transformers 、 datasets 、 evaluate 库实现预训练模型的加载、微调 and 评估。在 200 条数据上测试最终可以得到较高的正确率：

```
(bupt-ml 3.12.0)
# dsy @ DESKTOP-Q6TLA7T in ~\Code\bupt\bupt-ml on git:main [13:51:03]
$ python .\test.py
200it [00:27, 7.17it/s]
Accuracy: 1.0
```