

实验二 基于 chatglm3-6b 模型的 lora 方法微调

一、实验内容

ChatGLM3 是智谱 AI 和清华大学 KEG 实验室联合发布的对话预训练模型。ChatGLM3-6B 是 ChatGLM3 系列中的开源模型，在保留了前两代模型对话流畅、部署门槛低等众多优秀特性的基础上，ChatGLM3-6B 引入了更强大的基础模型；更完整的功能支持；更全面的开源序列，并且对学术研究完全开放。

LoRA (Low-Rank Adaptation of Large Language Models, 大型语言模型的低秩适应)是微软研究员提出的一种新颖技术,旨在解决微调大型语言模型的问题。LoRA 提议冻结预训练模型的权重,并在每个 Transformer 块中注入可训练层(称为秩分解矩阵)。这大大减少了可训练参数的数量和 GPU 内存需求,因为大部分模型权重不需要计算梯度。通过专注于大型语言模型的 Transformer 注意力块,LoRA 的微调质量与完整模型的微调相当,同时速度更快,计算需求更低。

本次实验参考 ChatGLM3 官方提供的代码以及微调手册,带领大家学习 Lora 微调的方法流程,为大家学习大语言模型打好基础。

参考文件:

ChatGLM3 官方地址: [THUDM/ChatGLM3: ChatGLM3 series: Open Bilingual Chat LLMs](https://github.com/THUDM/ChatGLM3)
| [开源双语对话语言模型 \(github.com\)](https://github.com/THUDM/ChatGLM3)

Lora 原文: <https://arxiv.org/abs/2106.09685>

Step1 打开终端, 下载代码

打开 **Terminal** 终端, 把 ChatGLM3.zip 文件上传到服务器, 并解压
unzip ChatGLM3.zip
cd ChatGLM3

Step2 创建虚拟环境并激活(2 分)

虚拟环境命名使用: 姓名首字母小写+学号, 如:cqf1190202318

conda create -n cqf1190202318 python=3.10

conda activate cqf1190202318

创建成功后给出截图。如:

```
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
#
# To activate this environment, use
#
#     $ conda activate cqf1190202318
#
# To deactivate an active environment, use
#
#     $ conda deactivate

root@dsw-346519-b567d758-fw7cz:/mnt/workspace/ChatGLM3# conda activate cqf1190202318
(cqf1190202318) root@dsw-346519-b567d758-fw7cz:/mnt/workspace/ChatGLM3#
```


Step6 安装微调所需的依赖

本实验使用官方提供的微调代码，按照如下指令安装依赖。

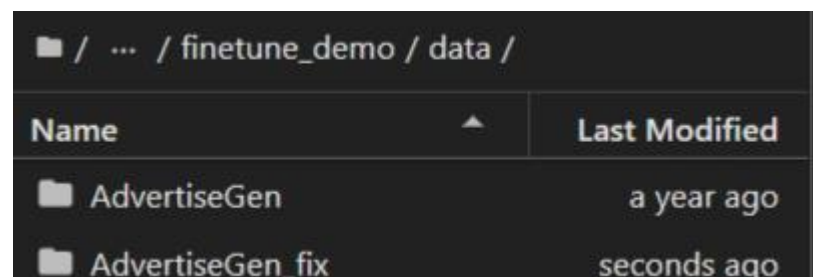
```
cd ../
cd finetune_demo/
conda install gcc_linux-64
pip install -r requirements.txt
pip install nltk # 补充安装的库，不然后续训练会报错
```

Step7 下载并处理 AdvertiseGen 数据集（2 分）

```
wget https://cloud.tsinghua.edu.cn/f/b3f119a008264b1cabd1/?dl=1
mv index.html?dl=1 AdvertiseGen.tar.gz
mkdir data
tar -xzf AdvertiseGen.tar.gz -C /mnt/workspace/ChatGLM3/finetune_demo/data
然后需要对解压后的 AdvertiseGen 数据集进行处理，以符合模型训练所需格式，
数据需转化到如下格式：
```

```
{"conversations": [{"role": "user", "content": "类型#裙*裙长#半身裙"}, {"role": "assistant", "content": "这款百搭时尚的仙女半身裙……。"}]}
```

数据可以自行处理，也可以按照官方提供的 `lora_fineyune.ipynb` 处理数据集，切割后的数据需保存在 `finetune_demo/data/AdvertiseGen_fix` 文件下。数据处理成功后，`data` 文件夹会出现以下两个文件。



报告部分需提供你处理数据所使用的代码截图和运行成功后的截图。

Step8 训练模型（4 分）

运行下面指令训练模型：

```
python finetune_hf.py 处理后的数据集路径 下载的 chatglm3-6b 模型路径 configs/lora.yaml
```

如果出现报错，按照提示安装所需函数库

模型成功训练需截图。

Step9 验证广告生成功能（4 分）

```
python inference_hf.py your_finetune_path --prompt your prompt
```

可以使用如下提示词验证模型：

```
--prompt "类型#裙*版型#显瘦*材质#网纱*风格#性感*裙型#百褶*裙下摆#压褶*
裙长#连衣裙*裙衣门襟#拉链*裙衣门襟#套头*裙款式#拼接*裙款式#拉链*裙款
式#木耳边*裙款式#抽褶*裙款式#不规则"
```

```
--prompt "类型#裤*材质#羊毛*裤长#九分裤*裤口#微喇裤"
```

广告成功生成需截图