

大模型微调身份证问答实践

廖家源

25111501

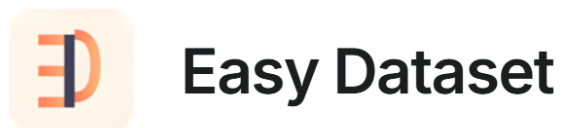
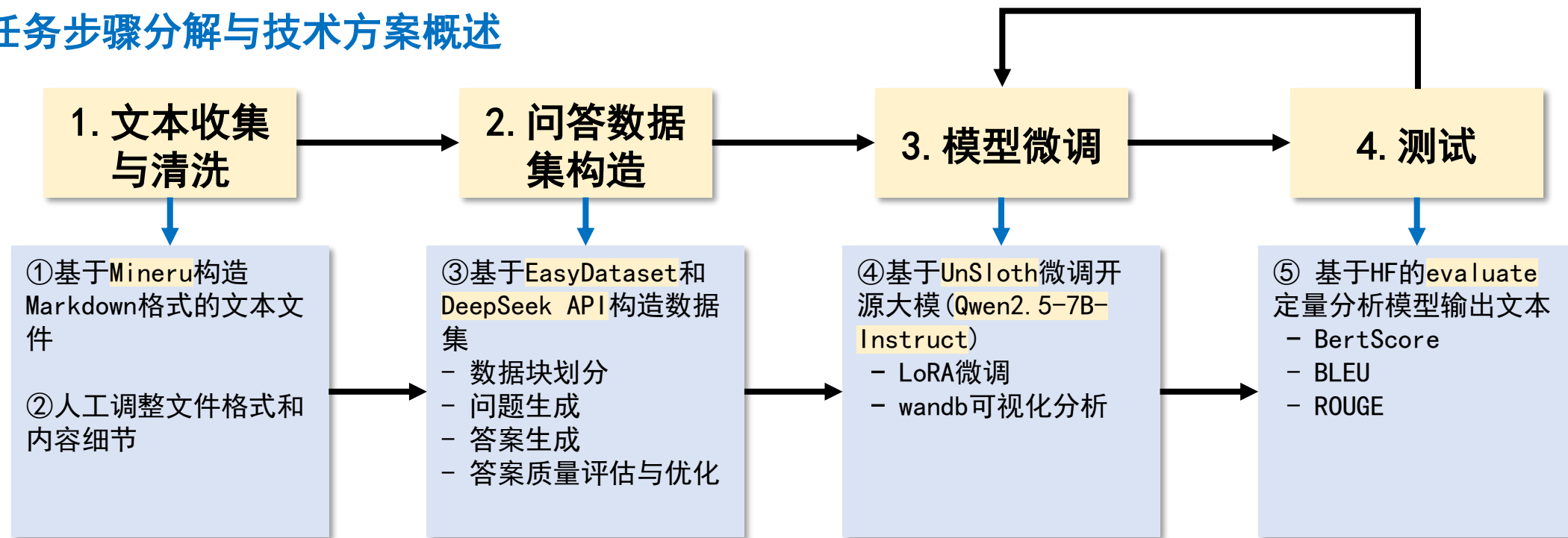
2025年12月26日

0.实现方案概述

□ 目标

- 提供居民身份证办理的相关文件数据，任务是用此数据对大模型进行微调，让大模型代替政府工作人员回答身份证业务相关问题。

□ 任务步骤分解与技术方案概述



1.文本收集与清洗

① 自动化格式转换

- ❑ 使用开源工具MinerU，将原文件转换为markdown格式，支持后续清洗和信息提取

② 手工清洗与格式化

- ❑ 删除署名备注、发布日期和目录等业务无关的文本
- ❑ 修正符号格式，删除空行等多余字符
- ❑ 上一步转换所得文件的标题均为一级标题，需要手动修正为多级标题形式，支持后续的基于格式层级的分块
- ❑ 其它细节上的小改动

关于在办理户口居民身份证业务工作中进一步严格规范冷僻字使用的通知.doc
中华人民共和国临时居民身份证管理办法.docx
现役军人和人民武装警察居民身份证申领发放办法.docx
关于转发广东省户口居民身份证管理工作操作规范（2024年版）的通知【正文】4433808.wps
关于全面上线居民身份证换补领业务和户籍三项“跨省通办”服务的通知（粤公网发〔2023〕1491号）.doc
关于开通首次申领居民身份证申领临时居民身份证“跨省通办”“全省通办”服务的通知（粤公网发〔2023〕654号）.doc



1.中华人民共和国临时居民身份证管理办法.md
2.关于全面上线居民身份证换补领业务和户籍三项“跨省通办”服务的通知（粤公网发〔2023〕1491号）.md
3.关于在办理户口居民身份证业务工作中进一步严格规范冷僻字使用的通知.md
4.关于开通首次申领居民身份证申领临时居民身份证“跨省通办”“全省通办”服务的通知（粤公网发〔2023〕654号）.md
5.现役军人和人民武装警察居民身份证申领发放办法.md
61.关于转发广东省户口居民身份证管理工作操作规范（2024年版）的通知.md
62.关于转发广东省户口居民身份证管理工作操作规范（2024年版）的通知.md

2.问答数据集构造

① 数据块划分与清洗

- ❑ 根据markdown文档结构，把篇幅较长的文本拆解成若干较小的、便于处理的文本块
- ❑ 利用大模型（DeepSeek-V3.2）清洗文本

② 问题生成

- ❑ 基于每个文本块，利用大模型（DeepSeek-V3.2）生成问题

③ 答案生成

- ❑ 基于每个文本块和上一步生成的问题，利用大模型（DeepSeek-R1）生成答案

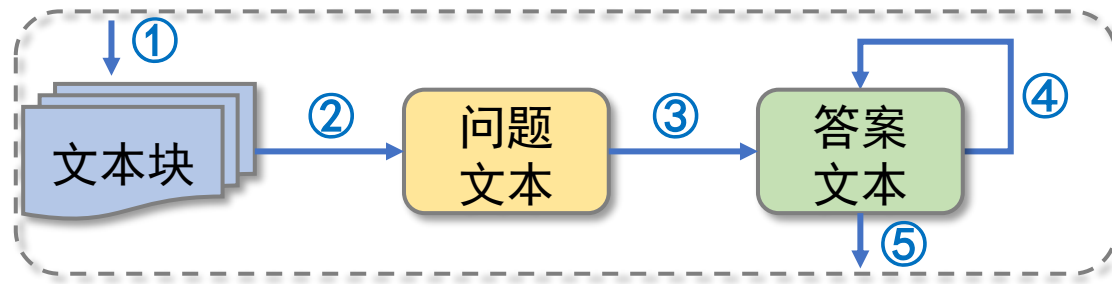
④ 答案质量评估与优化

- ❑ 利用大模型（DeepSeek-R1）评估每个问题-答案对，若原答案质量较差，则给出改进意见指导重新生成高质量答案。

⑤ 问答数据集构造

- ❑ 将问答对保存为Alpaca格式的数据集

1. 中华人民共和国居民身份证管理办法.md
2. 关于全面上线居民身份证换补领业务和户籍三项“跨省通办”服务的通知（粤公网发〔2023〕1491号）.md
3. 关于在办理户口居民身份证业务工作中进一步严格规范冷僻字使用的通知.md
4. 关于开通首次申领居民身份证申领临时居民身份证“跨省通办”“全省通办”服务的通知（粤公网发〔2023〕654号）.md
5. 现役军人和人民武装警察居民身份证申领发放办法.md
61. 关于转发广东省户口居民身份证管理工作操作规范（2024年版）的通知.md
62. 关于转发广东省户口居民身份证管理工作操作规范（2024年版）的通知.md



```
1 > {  
2 > {  
3     "instruction": "公民个人申请查询本人家庭户籍档案时，需要提交哪些材料？",  
4     "input": "",  
5     "output": "公民个人申请查询本人家庭户籍档案时，需要提交以下材料：\n1. 申办人的  
6     "system": ""  
7 }  
8 > { ...
```

（约2.2k个问答对）

 Easy Dataset <https://github.com/ConardLi/easy-dataset>

 deepseek [API: deepseek-ai/DeepSeek-V3.2 & deepseek-ai/DeepSeek-R1](https://api.deepseek-ai.com/DeepSeek-V3.2)

3.模型微调

□ 实验配置

- 训练框架: Unsloth (开源的大模型微调工具, 支持高效、低显存训练)
- 微调算法: LoRA
- 基座模型: Qwen2.5-7B-Instruct
- 实验环境: Ubuntu 22.04.5 LTS, NVIDIA A100 80GB, CUDA_12.6, PyTorch_2.9.0+cu128

□ 数据集处理

- 将上一步构造所得数据集随机划分为训练集和验证集两部分, 两者分别占比80%和20%
- User Prompt设置 (不修改模型的System Prompt)

```
alpaca_prompt_domain_special2 = \
"""你是一名专门负责居民身份证相关业务的政务服务助手。针对用户提问, 请遵循以下规范进行答复:
1. 解答内容必须基于居民身份证办理的官方政策, 不得编造或杜撰规定。
2. 回答应当简洁准确且清晰规范。
3. 如用户提问与居民身份证业务无关, 请礼貌说明并避免提供不确定的信息。
### 用户提问
{instruction}
### 答复
{output}
"""
```

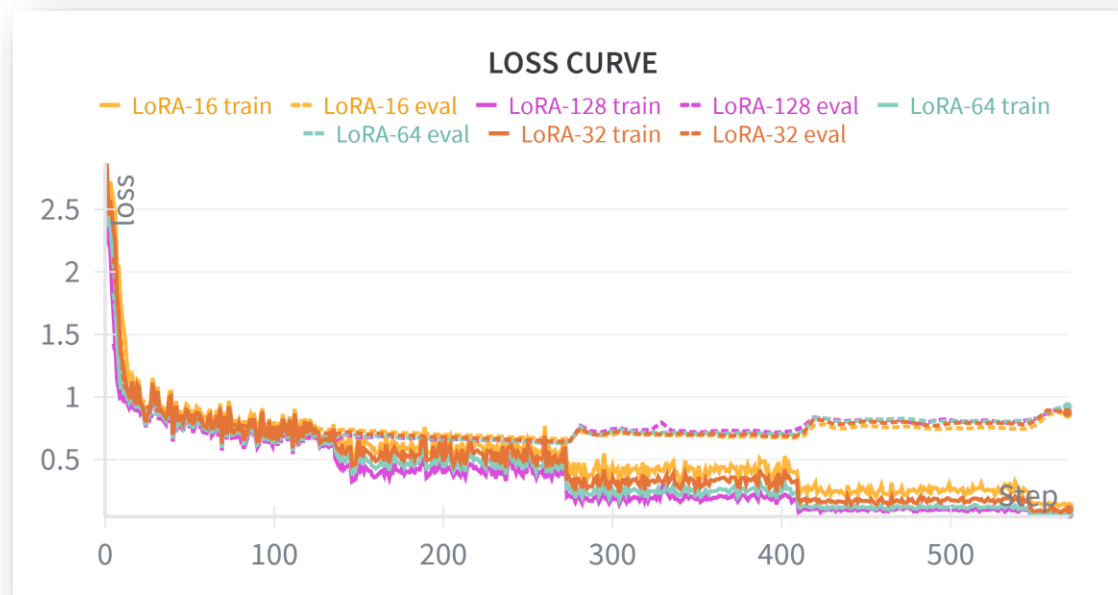


3.模型微调

□ 关键参数配置

- seed=32
- `lora_rank=16/32/64/128`; `lora_alpha=lora_rank*2`; `lora_dropout=0`; `bisa=None`
- `lora_target_modules= (all linear layers in Attention & FFN)`
- `lr=2e-4`; `lr_scheduler_type=Linear`
- `epoch=5 (114 step/epoch)`
- `batch_size=16=2*8`

□ 训练效果



- 四种不同LoRA Rank的配置的损失曲线走势接近
- 每个epoch有114次step, 经过近3个epoch的微调后, 开始呈现过拟合趋势

4.测试

□ 评测指标：BERT Score (F1/Precision/Recall)、BLEU、ROUGE-L

测试集

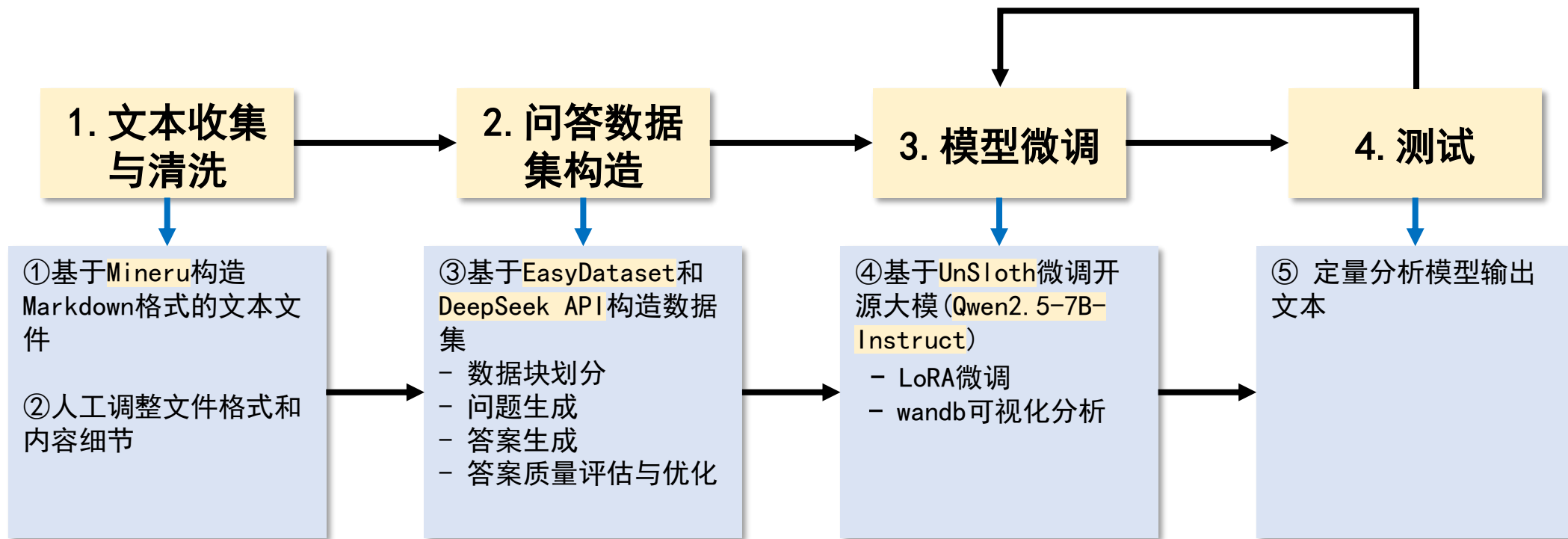
Model	F1	Precision	Recall	BLEU	ROUGE-L
Base	0. 7060	0.6797	0. 7404	0.0347	0.0818
LoRA-16	0. 7767	0.7798	0. 7810	0.1042	0.0893
LoRA-32	<u>0. 7772</u>	<u>0.7779</u>	0. 7839	<u>0.1217</u>	0.1089
LoRA-64	0. 7796	0. 7778	<u>0. 7892</u>	0.1285	0.1282
LoRA-128	0. 7717	0. 7486	0. 8023	0.1174	<u>0.1228</u>

完整数据集

Model	F1	Precision	Recall	BLEU	ROUGE-L
Base	0.7046	0.6781	0.7391	0.0358	0.0698
LoRA-16	0.8709	0.8794	0.8669	0.1841	0.1312
LoRA-32	0.9074	0.9082	0.9092	0.3912	0.1833
LoRA-64	<u>0.9231</u>	<u>0.9301</u>	<u>0.9189</u>	<u>0.4600</u>	<u>0.2027</u>
LoRA-128	0.9343	0.9358	0.9350	0.5640	0.2237

- BERTScore显著高于ROUGE/BLEU，说明模型主要提升了语义一致性，而非简单记忆或复现标准答案。
- LoRA微调在身份证业务问答任务上明显提升模型性能，在语义一致性（BERTScore）与文本重叠度（BLEU、ROUGE-L）指标上均优于基础模型。
- 在完整数据集上，模型性能随 LoRA Rank 增大呈现稳定提升趋势，表明更高 Rank 在数据充分时具有更强的领域建模能力。
- 相比完整数据集，模型在测试集上的性能提升幅度较小，提示其泛化能力仍有进一步优化空间。

5.总结



❑ 文本清洗与数据集构造是模型微调任务的主要工作量和难点所在

- ❑ Mineru的文件格式自动化转换会导致文本章节层级错乱，需要人工检查并调整，花费较多时间
- ❑ EasyDataset的数据块划分策略不完善，导致相同业务下不同小节的文本段被切分到不同文本块，语义不连贯
- ❑ 基于大模型生成的问题、答案质量参差不齐；基于质量评估的答案重写的优化效果不明显

❑ 模型微调部分相关超参数仍需调优，提升模型在测试集的性能表现

❑ 未使用RAG等方式提升模型能力