

ADL HW3 Report

NTU CSIE, R12922051
資工碩一 陳韋傑

Q1: LLM Tuning

Parameters

- Bits: 4
- Lora Rank: 64
- Learning Rate: 3e-5
- Source Max Length: 1024
- Target Max Length: 256
- Batch Size: 16
- Max Step: 1000 (2.01 Epoch)
- Optimizer: paged_adamw_32bit

I use QLoRA, a PEFT method, to tune Taiwan-LLaMA. For every linear layers, I add a LoRA layer and only finetune them.

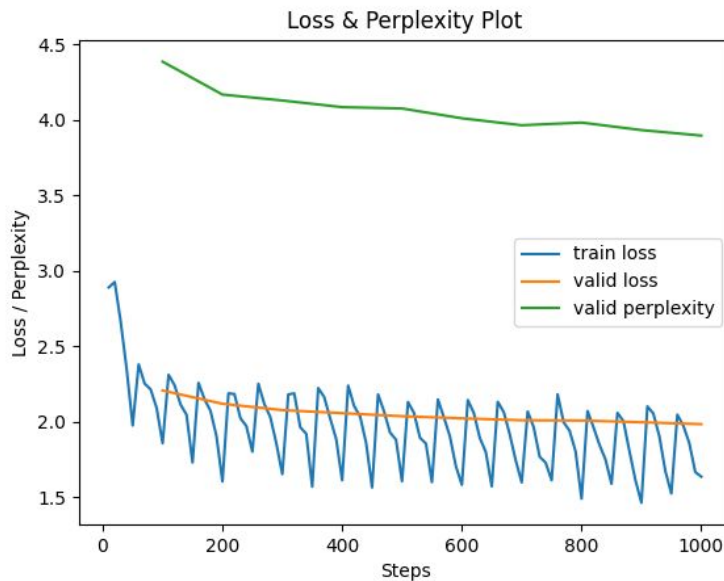
Q1: LLM Tuning

Total Steps: 1000

Epoch: 2.01

At Step 1000:

- Train Loss: 1.6367
- Valid Loss: 1.9845
- Perplexity: 3.8961



Q2: LLM Inference Strategies

Setting	Perplexity
Zero-shot on Taiwan-LLaMA	5.0744
Few-shot on Taiwan-LLaMA (Shot=1)	1.1323
Few-shot on Taiwan-LLaMA (Shot=3)	1.1756
Few-shot on Taiwan-LLaMA (Shot=5)	1.1775
QLoRA on Taiwan-LLaMA	3.8961

Q2: LLM Inference Strategies

Settings:

- Zero-Shot
 - Same prompt as QLoRA (Since there are no example provided)
- Few-Shot
 - Refined prompt for model understanding and examples
 - Different numbers of example are provided, example are chosen in the relative position in the dataset (previous data pair)

Comparison

Between three methods, Few-Shot in-context learning performs the best (no matter how many examples are used), while QLoRA performs worse than Few-Shot but better than Zero-Shot.

Q2: LLM Inference Strategies - Prompt

- Zero-Shot & QLoRA

你是一個精通現代中文與文言文的大師，以下是用戶和你之間的對話。你的目標是對用戶的問題提供有用、精確且簡潔的回答。USER: {instruction} ASSISTANT:

- Few-Shot

你是一個精通現代中文與文言文的大師，以下是用戶和你之間的對話。你的目標是對用戶的問題提供有用、精確且簡潔的回答。以下為幾個翻譯的正確例子
USER:{example}
ASSISTANT:{answer}\n) x N 以下為你需要翻譯的句子，請根據前面提供的正確結果進行翻譯
USER:{instruction} ASSISTANT:

Reference

- <https://github.com/MiuLab/Taiwan-LLaMa>
- <https://github.com/artidoro/qlora>