

ADL HW3

r12944005 陳乙馨

Q1 : LLM Tuning

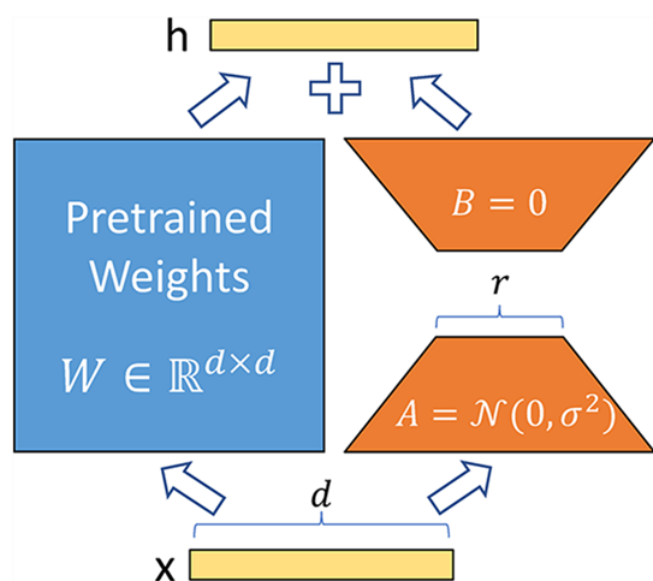
How much training data did you use?

我使用了全部 10000 筆 training data 來進行微調。

How did you tune your model?

我主要是使用 pretrained model `zake7749/gemma-2-2b-it-chinese-kyara-dpo`，並採用 QLoRA 方法進行微調。LoRA 是一種 LLM finetuning 的技術，他的優勢在於無須改變原本 pretrained model 的參數，也無須在 pretrained model 裡面增加 adapter layer 造成額外的延遲，LoRA 只需要額外訓練兩個較小的矩陣 A 和 B，使其相乘後和 pretrained model 的維度 W 相同即可，最後的輸出僅是將兩者加總 $W + AB = Y$ ，這樣的設計允許 A 和 B 與 pretrained model 的原始權重平行運行，有效減少計算負擔並提高效能。

QLoRA 則是進一步透過 Quantization 減少 memory usage。它支援像 4 bit 或 8 bit 的量化，使得訓練過程更加輕量，特別適合在資源受限的環境中部署大規模模型。



- 一開始我使用助教給的 default prompt，然而我發現 perplexity 在 `public_test.json` 上的表現到 25 就下不去了，因此我請 LLM 幫我生成了客製化的 prompt，並不斷調整 learning rate 之後，perplexity 最低才來到 14，以下是我修改後的 prompt：

```
"你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。"
\nUSER :
{instruction}\nASSISTANT : "
```

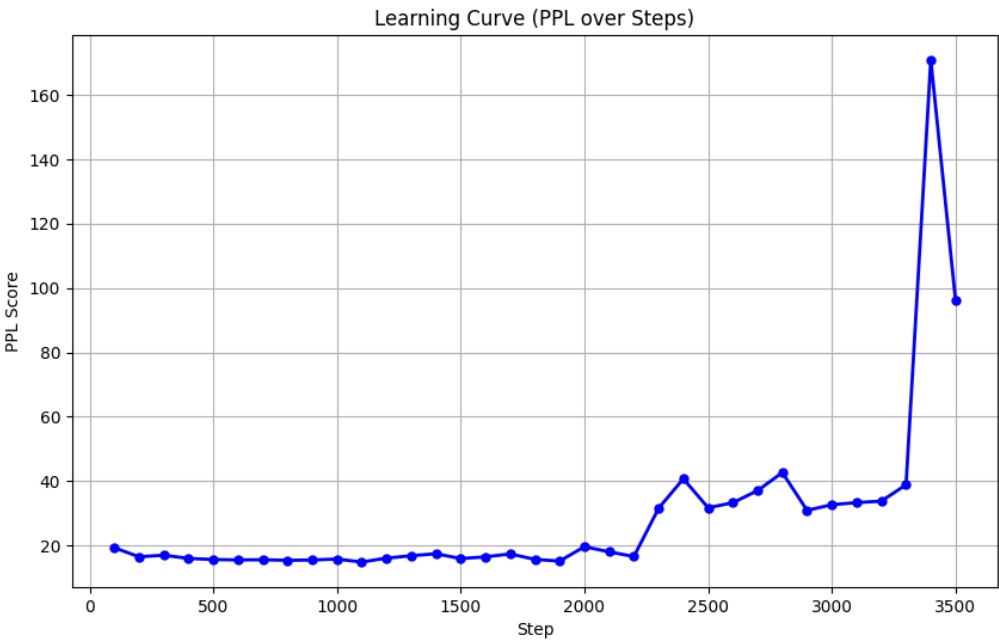
- 經過幾輪訓練，我發現在約 3000 steps 後 perplexity 就會開始指數性飆高，推測是因為 overfitting，因此我將 `--max_step` 設為 3500，以展示模型從慢慢收斂到後面開始 overfitting 的學習過程。

What hyper-parameters did you use?

Pretrained Model	zake7749/gemma-2-2b-it-chinese-kyara-dpo
lora_r	64
lora_alpha	16
lora_dropout	0.0
per_device_train_batch_size	4
gradient_accumulation_steps	2
max_steps	3500
learning_rate	1e-4
lr_scheduler	constant

What is the final performance of your model on the public testing set?

Best Checkpoint	Best Perplexity on public testing set
1100 steps	14.83434375



Q2: LLM Inference Strategies

What is your setting?

依然是使用 Pretrained Model : zake7749/gemma-2-2b-it-chinese-kyara-dpo · BitsAndBytes 設定如下 (和微調時的設定一致) :

```
quantization_config=BitsAndBytesConfig(
    load_in_4bit=True,
    load_in_8bit=False,
    bnb_4bit_compute_dtype=torch.bfloat16,
    bnb_4bit_use_double_quant=True,
    bnb_4bit_quant_type='nf4',
)
```

How many in-context examples are utilized? How did you select them?

- 在 Few-Shot Learning 的部分，我主要是從 `train.json` 裡面挑選範例，並且分別對 `#Example = 1, 2, 4` 進行測試。
- 在 `#Example = 2` 時，我特別選擇文言文轉白話文、以及白話文轉文言文兩種不同的範例。

Prompt Design / Comparison

- 從以下的結果可知，In-context learning 在 2B 小模型的表現並不理想，表現遠比 QLoRA 還要差。
- 進行 In-context learning 時，Examples 的數量確實會影響 Performance，由下表可以觀察出 Few-Shot 的結果比 Zero-Shot 好很多，且隨著 example 數量增加，perplexity 也有下降的趨勢。

Strategy	#Example	Prompt	Perplexity
QLoRA	0	你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。 無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。 USER : {instruction}\nASSISTANT :	14.83434375
Zero-Shot	0	你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。 無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。 USER : {instruction}\nASSISTANT :	789.9786875
Few-Shot	1	你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。 無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。 USER : 沒過十天，鮑泉果然被拘捕。 ASSISTANT : 後未旬，果見囚執。 USER : {instruction}\nASSISTANT :	328.3041875

Strategy	#Example	Prompt	Perplexity
Few-Shot	2	<p>你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。</p> <p>USER：沒過十天，鮑泉果然被拘捕。 ASSISTANT：後未旬，果見囚執。</p> <p>USER：辛未，命吳堅為左丞相兼樞密使，常林參知政事。 ASSISTANT：初五，命令吳堅為左丞相兼樞密使，常增為參知政事。</p> <p>USER：{instruction} ASSISTANT：</p>	296.973875
Few-Shot	4	<p>你是一位語言轉換助理，能夠根據用戶的輸入，自動判斷句子是白話文還是文言文、並且進行互譯。無論用戶提供的是白話文還是文言文，你都要進行正確的轉換並給出簡潔的翻譯。</p> <p>USER：沒過十天，鮑泉果然被拘捕。 ASSISTANT：後未旬，果見囚執。</p> <p>USER：辛未，命吳堅為左丞相兼樞密使，常林參知政事。 ASSISTANT：初五，命令吳堅為左丞相兼樞密使，常增為參知政事。</p> <p>USER：文言文翻譯： ASSISTANT：明日，趙用賢疏入。 ASSISTANT：第二天，趙用賢的疏奏上。</p> <p>USER：翻譯成現代文： ASSISTANT：渭州人鄭五醜造反，與叛逆羌傍乞鐵忽互相呼應。下令趙剛前往鎮壓。</p> <p>USER：{instruction} ASSISTANT：</p>	260.6014375

Q3 : Try Llama3-Taiwan

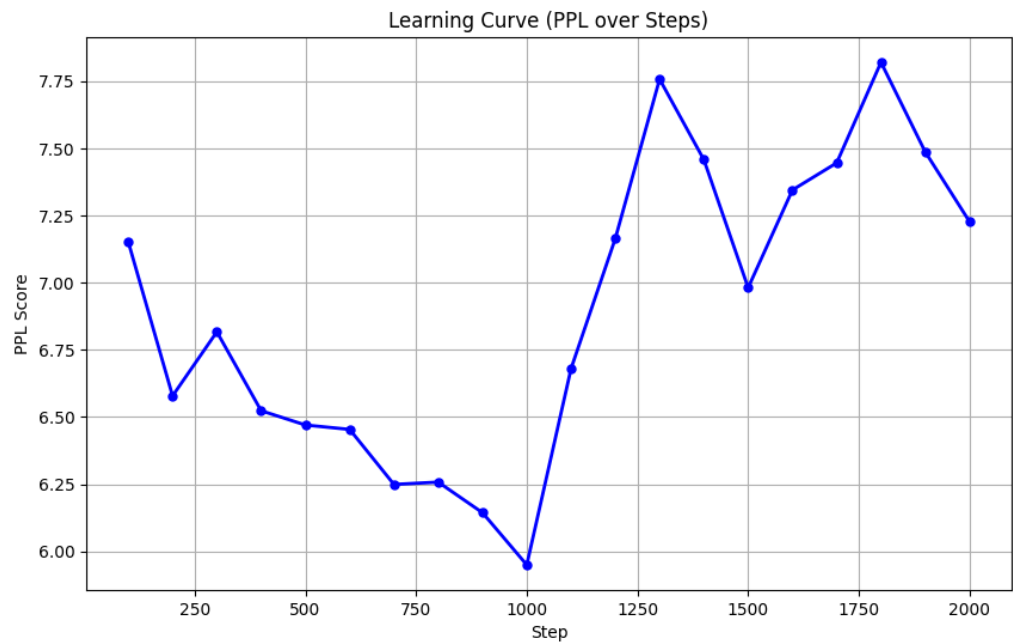
在這個部分，我使用了 `yentinglin/Llama-3-Taiwan-8B-Instruct` 進行 QLoRA 的 finetune，由於該 model 較大、運算資源不足的緣故，只能選用較小的 `batch size` 跟 `lora rank`。

Hyper-parameters

Pretrained Model	yentinglin/Llama-3-Taiwan-8B-Instruct
<code>lora_r</code>	16

Pretrained Model	yentinglin/Llama-3-Taiwan-8B-Instruct
lora_alpha	16
lora_dropout	0.0
per_device_train_batch_size	1
gradient_accumulation_steps	8
max_steps	2000
learning_rate	1e-4
lr_scheduler	constant

Learning Curve



Best Checkpoint	Best Perplexity on public testing set
1000 steps	5.94946875

- 此模型表現最好的 checkpoint 也在 1000 steps 左右，再繼續訓練下去又會有 overfitting 的狀況。
- 他的訓練成果明顯優於 `gemma-2-2b-it-chinese-kyara-dpo`，並且在使用 `public_test.json` 進行 inference 後，可以發現這個模型的回答內容可讀性通常較高，比較不會出現看不懂的字詞。以下比較兩個 model 在同一筆 public test data 的作答狀況：

```
{
  "id": "87945a20-f869-4be9-b586-f5ce20ddd78b",
  "instruction": "文言文翻譯：\n中宗與庶人嘗因正月十五日夜幸其第，賜賚不可勝數。",
  "output": "答案：唐中宗與韋庶人曾經在正月十五日夜到韋安石的宅第，並賜賞給他不可勝數的財物。"
}
```

Model	<code>zake7749/gemma-2-2b-it-chinese-kyara-dpo</code>	<code>yentinglin/Llama-3-Taiwan-8B-Instruct</code>
Answer	答案：憲宗與庶民曾經在大年初一，晚上嚮他住的府裏拜訪，皇帝賞賜的文禮物和禮物無法列舉。	唐中宗曾在正月十五日夜裏去過他府上，賜給他的東西不可勝數。