

* Week 2! - Fine tuning LLMs with instruction

* Instruction fine-tuning

- Limitations of In-context Learning
(One shot or few shot inference)

- May not work for smaller LLM
- Examples may take up space in context window.

- Fine-tuning is supervised learning technique where you use set of labeled examples to update the weights of the LLM.

- The labeled examples are prompt completion pair.

• Instruction fine-tuning

- Instruction fine-tuning trains the model using examples that demonstrate how it should respond to specific instruction.

for ex.

Instruction → classify this review:

I loved this DVD.

Desired caption → Sentiment: Positive.

- Process like Instruction fine-tuning where all the weights of model gets updated are known as Full fine-tuning.

- LLM fine-tuning process
 - Prepare Instruction dataset
 - Divide the dataset into Training, Validation & Test splits.
 - During fine-tuning pass the prompt to pretrained LLM you will get completion. Now compare this with the desired label in training dataset.

- calculate loss using standard cross-entropy
 - update model ~~as~~ weights from the calculated loss in standard backpropagation.
- * Fine tuning on a single task
- Fine tuning for a single task can be beneficial and also can be achieved with very few examples
 - But this ~~is~~ can cause a phenomenon called catastrophic forgetting.
- Catastrophic forgetting
 - Fine tuning can significantly increase the performance of a model on a specific task.
 - but can lead to reduction in ability on other task.

• How to avoid catastrophic forgetting

- First note that you may not have to.
- Fine tune multiple tasks at the same time.
- Consider Parameter Efficient Fine tuning (PEFT)

* Multi-task instruction fine tuning

- multi-task finetuning is an extension of single task fine tuning where training data is comprised of example input & outputs for multiple tasks.
- one ~~drawback~~ drawback of multi task fine tuning is you will require a lot of data. You may need as many as 50-100,000 examples in your training set.
- For ex. FLAN family of models are trained using multi-task fine tuning.

* Model Evaluation.

- LLM Evaluation - challenges

- Since LLM works on Natural languages it is hard to evaluate LLMs.

- LLM Evaluation - metrics

ROGUE

- Used for text summarization
- Compare a summary to one or more reference summaries

BLEU SCORE

- Used for text translation
- Compares to human generated translations.

unigram → one word
 bi gram → two words
 n-gram → group of n-words

$$\text{ROGUE 1} = \frac{\text{unigram matches}}{\text{unigrams in reference}}$$

Recall

$$\text{ROGUE 1} = \frac{\text{unigram matches}}{\text{unigrams in output}}$$

Precision

$$\text{ROGUE 1} = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

F₁

- Similar to ROGUE, we can calculate ROGUE₂ with help of bi-grams
- ROGUE-L, where L is Longest common Subsequence (LCS)

{ Longest common subsequence (LCS) is
 the longest sequence of words that appear in both sentences in same order, but not necessarily, consecutive }

- ROGUE clipping.

BLEU SCORE.

- The BLEU SCORE quantifies the quality of translation by checking how many n-grams in machine-generated translation match those in reference translation
- $\text{BLEU}(\text{metric}) = \text{Avg}(\text{precision across range of metric})_{\text{n-gram sizes}}$

* Benchmarks

- In order to evaluate LLM's more holistically, we can make use of pre-existing datasets & associated benchmarks.
- Benchmarks such as GLUE, SuperGLUE & HELM, cover a wide range of task scenarios.

- Benchmarks for massive models.
 - massive multitask Language Understanding (MMLU) 2021
 - Big-bench, Big-bench hard, Big-bench lite 2022
 - ~~Holistic~~

Holistic Evaluation of Language models (HELM)

Metrics

- Accuracy
- Calibration
- Robustness
- Fairness
- Bias
- Toxicity
- Efficiency

* Parameter Efficient fine-tuning

- PEFT only updates small set of parameters
- In some cases only 15-20% of original LLM weights are trained & rest are frozen.
- Less prone to Catastrophic forgetting.
- PEFT fine tuning saves space & is flexible.
- PEFT Tradeoffs
 - Parameter Efficiency
 - Memory Efficiency
 - model performance.
 - Inference Costs
 - Training Speed.

• PEFT methods

i) Selective methods

- Selective methods are those that fine tune only a subset of the original LLM parameters.
- There are several approaches that you can take to identify which parameters you want to update.
- You have option to train only certain components of the model or specific layers or even individual parameter types.
- Researchers have found that performance of this methods is mixed & there are significant tradeoffs between parameter efficiency & compute efficiency.

ii) Reparameterization Methods

- Reparameterization methods also work with original LLM parameters, but reduce the number of parameters to train by creating new low rank transformations of the original network weights.
- A commonly used technique of this type is LORA

iii) Additive Methods

- Additive methods carry out fine tuning by freezing all the original LLM weights & introducing new trainable components.
- Here there are two main approaches
 - Adapters
 - Self prompts

- Adapters → Add new trainable layers to the architecture of the model, typically inside the encoder or decoder components after the attention or feed-forward layers.
- Self prompters → Keeps the model architecture fixed or frozen & focus on manipulating input to achieve better result.
This can be done by adding trainable parameters to prompt embedding or keeping the input fixed & retraining embeddings weights.

* PEFT Techniques I: LORA

LORA \rightarrow Low Rank Adaptation of Large Language Models

- Freeze most of the original LLM weights
- Inject 2 rank decomposition matrix.
- Train the weights of the smaller matrices

Steps to update model for inference

- Matrix multiply the low rank matrices

$$B * \bar{A} = \boxed{B \times A}$$

- Add to original weights

Frozen
Weights

+ $\boxed{B \times A}$

For example if you want to fine-tune

Use the base transformer model presented in "Attention is all you need" paper

- Transformer weights have dimensions $d_{xk} = 512 \times 64$
- So $512 \times 64 = 32,768$ trainable parameters

In LORA rank $r = 8$

- A has dimensions $r \times k = 8 \times 64 = 512$
- B has dimensions $d_{xr} = 512 \times 8 = 4,096$ trainable parameters
- 86% reduction in parameters to train.

you can use LORA to train for many tasks. switch out the weights when you need them & avoid storing multiple full sized versions of LLM.

How to choose the rank of LORA

→ this is still a area of research.
In theory lower the rank, the smaller the number of trainable parameters.

* PEFT Techniques 2: soft prompts

- Prompt tuning adds trainable "soft prompt" to inputs.
- weights of model & soft prompt trained
- Instead embedding vectors of model gets updated.
- Very parameter Efficient strategy.
- performance of prompt tuning
 - prompt tuning doesn't perform as well as ~~as~~ full fine-tuning for smaller LMs.
 - But as model size reaches 10B param prompt tuning can be as efficient as full fine-tuning.