

СРАВНИТЕЛЬНЫЙ АНАЛИЗ МЕТОДОВ PEFT ДЛЯ ДООБУЧЕНИЯ БОЛЬШИХ ЯЗЫКОВЫХ МОДЕЛЕЙ

Маракулин А.А. (ИТМО), Дедкова А.В. (ИТМО), Аминов Н.С. (ИТМО)

Научный руководитель – Федоров Д.А.

(ИТМО)

Введение. На текущий момент созданы различные большие языковые модели (Large language models, LLM). Они показывают хорошие результаты на общем наборе задач, но, когда речь идет о запросах, выходящих за пределы стандартных сценариев, могут быть менее впечатляющими. Также может возникнуть необходимость добавить в модель новые данные. Эти проблемы можно решить дообучением имеющейся LLM. Однако обычный метод Fine-Tuning, при котором изменяются все веса, из-за размеров современных моделей становится слишком ресурсозатратным [1].

Методы PEFT (Parameter-Efficient Fine-Tuning) предоставляют эффективные решения для улучшения ответов LLM с использованием значительно меньшего объема ресурсов [1].

Основная часть. Для дообучения LLM существует множество методов. В общем случае мы можем просто продолжить обучение на новых данных и получить новую модель, однако такой подход упирается в имеющиеся вычислительные способности, так как многие модели имеют миллиарды параметров. В таком случае можно рассмотреть методы PEFT. LoRA (Low-Rank Adaptation), Prefix tuning, Prompt tuning, Adapters — основные из них, разберем их поподробнее.

Идея метода LoRA заключается в обучении дополнительной матрицы весов, представляя ее в виде произведения двух матриц меньшей размерности. При этом веса основной модели замораживаются [2]. Prefix tuning — аддитивный метод (аддитивные методы подразумевают под собой дополнение параметров модели новыми, обучение происходит именно на них, при этом исходные веса заморожены), в котором к каждому блоку трансформера мы добавляем последовательность обучающих векторов, называемых префиксами, и обучаем только их [3]. У Prefix tuning есть его упрощенная версия под названием Prompt tuning. Суть этого метода в добавлении к входным данным тензора подсказок (soft prompts, промпты сгенерированные ИИ), который и будет обучаться [4]. Adapters — еще один аддитивный метод обучения. В данном методе слегка изменяется архитектура трансформера: перед двумя блоками суммирования в кодировщике добавляется слой адаптеров. В итоге получается два новых слоя, первый переводит входной сигнал из измерения d в меньшее измерение m , а второй наоборот, из m в d [5].

Для дообучения были выбраны методы LoRA и Adapters ввиду отсутствия необходимости модификации входных данных. Обучение проходило на платформе Kaggle.

Выводы. В ходе исследования были проанализированы возможные методы дообучения PEFT. Из них были выбраны два метода и с их помощью обучена LLM, а также проведена оценка качества полученных моделей.

Список использованных источников:

1. Lialin V., Deshpande V., Rumshisky A. Scaling down to scale up: A guide to parameter-efficient fine-tuning //arXiv preprint arXiv:2303.15647. – 2023.
2. Hu E. J. et al. Lora: Low-rank adaptation of large language models //arXiv preprint arXiv:2106.09685. – 2021.
3. Lester B., Al-Rfou R., Constant N. The power of scale for parameter-efficient prompt

tuning //arXiv preprint arXiv:2104.08691. – 2021.

4. Li X. L., Liang P. Prefix-tuning: Optimizing continuous prompts for generation //arXiv preprint arXiv:2101.00190. – 2021.

5. Houshy N. et al. Parameter-efficient transfer learning for NLP //International Conference on Machine Learning. – PMLR, 2019. – C. 2790-2799.