

進捗報告

1 やったこと

- モデルマージ動作確認
- 小規模なパラメータ数の LLM のファインチューニング (AXCCEPT/llm-jp-3-3.7b-instruct-EZO-Humanities¹, gemma-2-2b-jpn-it²) のファインチューニング

2 モデルマージ動作

先週, 進化的モデルマージを試している記事³を参考に, モデルマージを試したが記事内の 7b のモデルを 3 つマージする構成では最後のマージ後のモデルをロードする部分で Cuda out of memory となった.

今週はとりあえず動かすことを最優先として,

- llm-jp-1.3b-v1.0⁴ を 2 つマージ
- 参考にした記事の 7b のモデルを 2 つに減らし, merge.method を dare.ties から linear に変える

上記の 2 つを試した.

結果として, 両方とも RTX3090 のサーバーで動いた. 図 1 に示すように後者の実験では学習時にタスクにおける score の上昇も見られ, 進化計算によりタスクに適したマージモデルを得られていることがわかる.

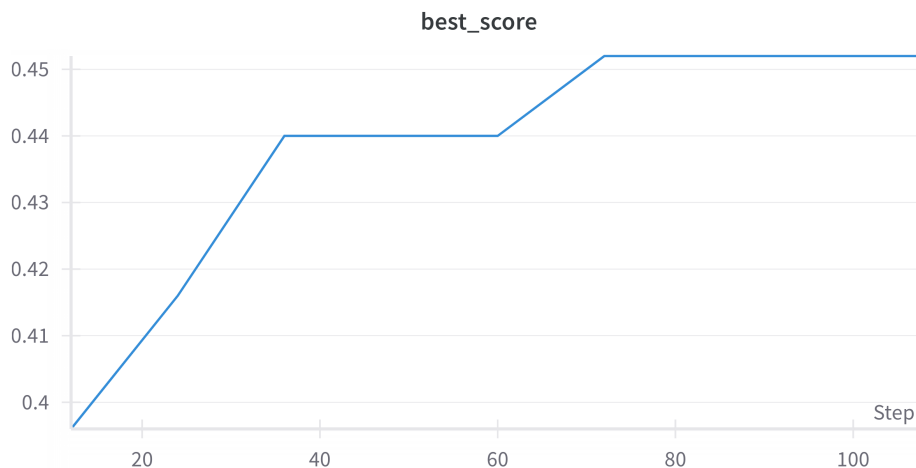


図 1: 学習時のタスクの score

なぜ GPU のメモリが 24 GB あれば動くはずの設定が動かなかったのか疑問が残るが, とりあえず動く設定が見つかったので, これを基に独自タスクで試していく予定です.

¹<https://huggingface.co/AXCCEPT/llm-jp-3-3.7b-instruct-EZO-Humanities>

²<https://huggingface.co/google/gemma-2-2b-jpn-it>

³<https://dalab.jp/archives/journal/llm-merge-evolve/>

⁴<https://huggingface.co/llm-jp/llm-jp-1.3b-v1.0>

3 小規模なモデルのファインチューニング

llm-jp-1.3b-v1.0 のモデルマージを試している際に、7b でのマージができない可能性を考慮し、7b 未満でのモデルのファインチューニングを試していた。試したモデルは以下の 2 つです。

- AXCEPT/llm-jp-3-3.7b-instruct-EZO-Humanities ⁵
- gemma-2-2b-jpn-it⁶

gemma-2-2b-jpn-it に関しては前回の報告で試し、出力が安定しない結果となった。今週は新たに以下の工夫をした。

- 学習データの LLM の出力部分に明示的に EOS トークンを追加
- retriever がクエリと関係ない話題を持ってくるのを防ぐように context を取ってくる方法を変更

校舎に関しては、これまではクエリと類似度が高いドキュメント上位 3 件を取ってくる形にしていたが、類似度のスコアを計算し一定の閾値を超えたものだけを取ってくる形にしてクエリと関係ないドキュメントを取ってこないようにした。

結果として、上記の工夫をしても出力に大きな変化は見られなかった。retriever の閾値を 0.8 としたことで関係ないドキュメントを引いてくることは少なくなったが、「あなたの名前はなんですか？」というクエリに対して、データベースには「名前は小野寺絢音」というドキュメントがあるにもかかわらずそのドキュメントを引いていなかった。EOS トークンの追加でも出力が安定せず、閾値も見直す必要がある。

また、llm-jp-3-3.7b-instruct-EZO-Humanities で retriever の閾値を 0.7 とした場合においても、出力が一文で完結しないケースが多く見られた。

AXCEPT-EZO-Common-9B-gemma-2-it ⁷ の学習では自然な応答ができていたので、モデルのパラメータ数による性能差もあるかもしれないが、ここまで同じ傾向が見られるということは学習時になにか問題が発生している可能性が高いため、学習時のコードを見直す必要がある。

4 今後やること

- モデルマージ (最優先)

パラメータ数 7b の 2 つのモデルであればマージできることがわかったため、7b でファインチューニングして出力が安定するモデルを見つける。また定量評価用の独自タスクを実装する必要があるためその実装。またファインチューニングのコードを見直す。

- RAG の Cross Attention のやつの実装

後回しにしてしまった。今週ファインチューニングがうまくいけばモデルマージの実験を長く待つだけになるので急いでやります。

参考文献

⁵<https://huggingface.co/AXCEPT/llm-jp-3-3.7b-instruct-EZO-Humanities>

⁶<https://huggingface.co/google/gemma-2-2b-jpn-it>

⁷<https://huggingface.co/AXCEPT/EZO-Common-9B-gemma-2-it>