

進捗報告

1 やったこと

- ファインチューニングのコードデバッグ
- モデルマージ

2 ファインチューニングのコードデバッグ

先週, 同じ Llama 2 派生の elyza/ELYZA-japanese-Llama-2-7b-instruct, stabilityai/japanese-stablelm-instruct-beta-7b を比較した結果, 前者はキャラクターを模したチャットボットとして自然な出力をしていたが, stabilityai/japanese-stablelm-instruct-beta-7b は出力の後ろに `i/s` がついているという結果になった. 学習の際のコードに問題があると考え, 詳しく調査していた.

2.1 損失の計算

SFT の際には LLM の応答部分だけ損失を計算しているようにしている. この際用いているのが, trl のライブラリの `DataCollatorForCompletionOnlyLM` を用いて学習データのバッチ処理をする際に, LLM の応答する部分以外に対して -100 のラベルをつけることで損失を計算する際にそれらの部分を無視できるようにしている.

学習データ

`i/s` `[INST]` `iiSYS` あなたは役立つアシスタントです. `iiSYS`
お嬢様のように振る舞ってほしいです. お嬢様が使うようなトーン、方式、語彙を使ってお嬢様のように
応答してほしいです. 応答の長さは 1 文程度で、30 文字程度に簡潔に回答してください. また必要に応じて
応答の参考になりうるお嬢様の情報を与えます. 応答の参考にならない情報を含む場合もあるので
その場合は無視してください.
こんにちは `[/INST]` こんにちは. お会いできて嬉しいですわ. `i/s`

[illegible]

損失の計算は, `trl` のライブラリにある `SFTTrainer` というクラスの継承元のクラスである `transformers` のライブラリの中の `Trainer` というクラスの `compute_loss` メソッドで処理されており¹, 具体的な Loss は `transformers` のライブラリの `LabelSmoother` というクラスが処理している。²

処理を見ると、モデルの出力の logits を負の対数尤度に変換し、ラベルで -100 となっているインデックスを避けながら、CrossEntropyLoss を計算し、損失を計算する部分で平均化することで Loss を計算している。

2.2 実際の Loss の推移

図 1, 2 に 2 つの LLM の SFT の際の Loss の推移を示す. 値は差があるものの, どちらも順調に減少しており大きな差は見られなかった.

3 今後やること

モデルマージの結果がまだ出ていない, 来週発表練習の予定なので結果が出てきたらまた個別に報告しに行きます。

参考文献

¹<https://github.com/huggingface/transformers/blob/v4.46.2/src/transformers/trainer.py#L3649>

²https://huggingface.co/transformers/v4.4.2/modules/transformers/trainer_pt_utils.html

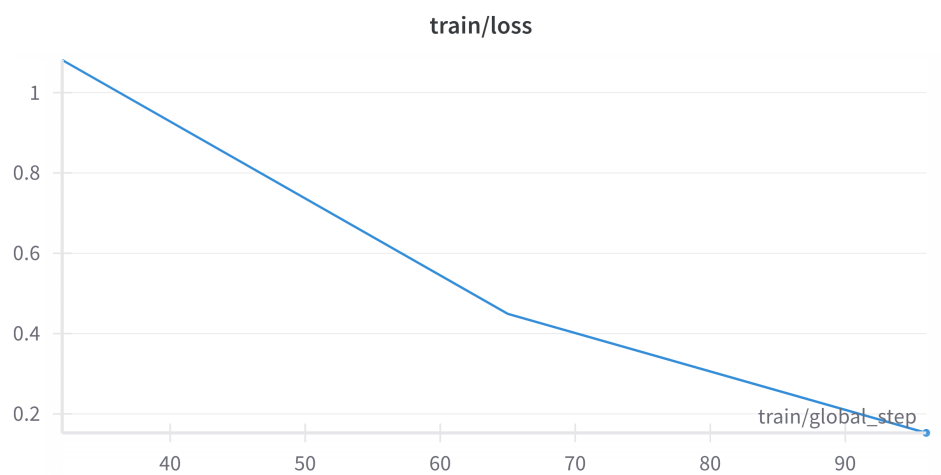


図 1: stabilityai/japanese-stablelm-instruct-beta-7b の Loss の推移



図 2: elyza/ELYZA-japanese-Llama-2-7b-instruct の Loss の推移