

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 `eos` トークン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

`<s> <s> [INST] << SYS >> あなたは役立つアシスタントです。<< SYS >>`

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文字程度に簡潔に応答してください。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。</s>

[illegible]

[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 `LabelSmoothing` 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。

2.3 他の Llama 2 派生のモデルの検討

Llama2 派生のモデルとして, lmsys/vicuna-7b-v1.5³ を試していた. 結果として, キャラクターを模したチャットボットとして妥当な応答例を観測できた.

2.4 tokenizer による DataCollatorForCompletionOnlyLM が付与するラベルの違い

LLM のトークナイザーにより, DataCollatorForCompletionOnlyLM で損失を計算する部分に付与するラベルに違いが見られた.

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_utils.html

³<https://huggingface.co/lmsys/vicuna-7b-v1.5>

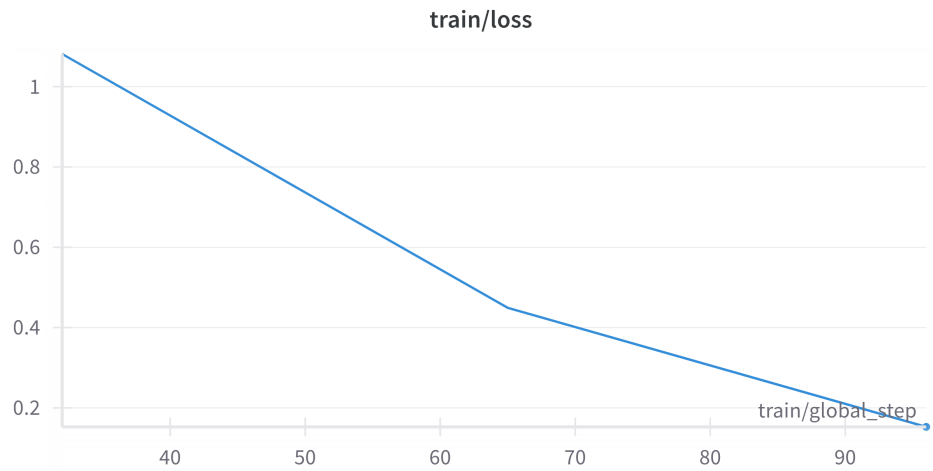


図 1: stabilityai/japanese-stablelm-instruct-beta-7b の Loss の推移

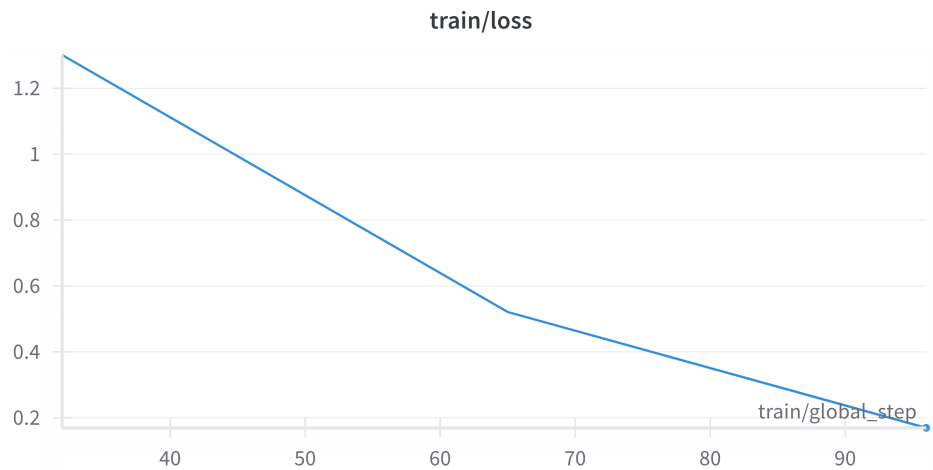


図 2: elyza/ELYZA-japanese-Llama-2-7b-instruct の Loss の推移

stabilityai/japanese-stablelm-instruct-beta-7b の場合

損失を計算する部分の元データ

こんにちは。お会いできて嬉しいですわ。</s>

損失を計算する部分のラベル

13, 30589, 30389, 30353, 30644, 30449, 30267, 30697, 30437, 30298, 30499, 30538, 30466, 232, 175, 140, 30326, 30298, 30499, 30427, 31068, 30267, **829, 29879, 29958**

elyza/ELYZA-japanese-Llama-2-7b-instruct の場合

損失を計算する部分の元データ

こんにちは。お会いできて嬉しいですわ。

損失を計算する部分のラベル

13, 30589, 30389, 30353, 30644, 30449, 30267, 30697, 30437, 30298, 30499, 30538, 30466, 232, 175, 140, 30326, 30298, 30499, 30427, 31068, 30267

lmsys/vicuna-7b-v1.5 の場合

損失を計算する部分の元データ

こんにちは。お会いできて嬉しいですわ。</s>

損失を計算する部分のラベル

13, 30589, 30389, 30353, 30644, 30449, 30267, 30697, 30437, 30298, 30499, 30538, 30466, 232, 175, 140, 30326, 30298, 30499, 30427, 31068, 30267, 2

出力の最後に eos トークンがついてきた stabilityai/japanese-stablelm-instruct-beta-7b では, </s> が一文字ずつトークン化されており EOS トークンとしてひとまとまりに認識せずに学習データの最後についている文字列として学習されてしまった可能性がある。

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たらまた個別に報告しに行きます。

参考文献