

# Direct Preference Optimization

Tri Nguyen

Oregon State University

February 7, 2024

# Alignment problem

- ▶ Human preference of a response  $\mathbf{y}$  given a prompt  $\mathbf{x}$  is measured by  $r_{\phi^{\natural}}(\mathbf{x}, \mathbf{y}) \geq 0$ .
  - ▶  $r(\mathbf{x}, \mathbf{y}_1) > r(\mathbf{x}, \mathbf{y}_2)$  means  $\mathbf{y}_1$  is more preferred than  $\mathbf{y}_2$ .
- ▶ Objective: given a trained language model  $\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})$ , fine-tune it so that
  - ▶ The outputs are aligned with human preference, while
  - ▶ Retaining the original model's generation skill.

A realized objective function:

$$\underset{\theta}{\text{maximize}} \quad \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [r_{\phi^{\natural}}(\mathbf{x}, \mathbf{y})] - \beta \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_{\text{kl}}(\pi_{\theta}(\mathbf{y} | \mathbf{x}) || \pi^{\text{ref}}(\mathbf{y} | \mathbf{x}))] \quad (1)$$

## Issues

1.  $r_{\phi^{\natural}}(\mathbf{x}, \mathbf{y})$  is unknown.
2. Problem (1) is “hard” to optimize due to the involvement of  $\theta$  in  $\mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})$  under expectation.

## The RL from Human Feedback approach [Ziegler et al. 2019]

- ▶ Estimate the score function  $r_{\phi^*}(\mathbf{x}, \mathbf{y})$
- ▶ Finetune the LLM model by optimizing the original objective function using the learned  $r_{\phi^*}$ .

## Fixing Issue 1: Specifying Preference Model

In hope of learning  $r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y})$ , we have to specify some model, and then obtain some samples.  
Preference Bradley-Terry model:

- ▶ Given  $L$  items, item  $i$  has a score  $s_i > 0$ .
- ▶ It models a binary result of an event  $i$  beats  $j$  as a Bernoulli RV with parameter

$$\Pr(i \succ j) = \frac{s_i}{s_i + s_j}, \quad \forall i, j \in [L].$$

In our LLM context,

$$\Pr(\mathbf{y}_1 \succ \mathbf{y}_2 \mid \mathbf{x}) = \frac{\exp(r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y}_1))}{\exp(r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y}_1)) + \exp(r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y}_2))} = \sigma(r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y}_2) - r_{\phi^\ddagger}(\mathbf{x}, \mathbf{y}_1)).$$

Under this model, the MLE objective is [\[Ziegler et al. 2019\]](#)

$$\underset{\phi}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mathcal{D}} \left[ I[\mathbf{y}_1 \succ \mathbf{y}_2] \sigma(r_{\phi}(\mathbf{x}, \mathbf{y}_2) - r_{\phi}(\mathbf{x}, \mathbf{y}_1)) + I[\mathbf{y}_2 \succ \mathbf{y}_1] \sigma(r_{\phi}(\mathbf{x}, \mathbf{y}_2) - r_{\phi}(\mathbf{x}, \mathbf{y}_1)) \right],$$

But there is no guarantee of learning the true  $r_{\phi^\ddagger}$ .

## Fixing Issue 2:

Now we have learned  $r_{\phi^*}$ , the objective is

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [r_{\phi^*}(\mathbf{x}, \mathbf{y})] - \beta \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_{\text{kl}}(\pi_{\theta}(\mathbf{y} | \mathbf{x}) \| \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ r_{\phi^*}(\mathbf{x}, \mathbf{y}) - \beta (\log(\pi_{\theta}(\mathbf{y} | \mathbf{x})) - \log(\pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))) \right] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [f_{\theta}(\mathbf{x}, \mathbf{y})]. \end{aligned}$$

This is a standard objective used in RL (policy gradient), hence can be solved using off-the-shelf tools such as PPO.

## A new approach

Rafael Rafailov et al. “Direct preference optimization: Your language model is secretly a reward model”. In: *arXiv preprint arXiv:2305.18290* [2023]

$$\underset{\pi_{\theta}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_{\text{kl}}(\pi_{\theta}(\mathbf{y} | \mathbf{x}) \| \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))] \quad (2)$$

This problem has “closed-form” solution:

$$\pi^*(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp\left(\frac{1}{\beta} r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})\right)$$

Note that RL people already known this, but this result is not very helpful due to the intractability of  $Z(\mathbf{x})$ .

## Proof of optimal policy

$$\begin{aligned}\arg \max_{\pi_{\theta}} \text{Objective} &= \arg \max_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} [r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})] - \beta \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [D_{\text{kl}}(\pi_{\theta}(\mathbf{y} | \mathbf{x}) \parallel \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))] \\ &= \arg \max_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y}) - \beta \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right] \\ &= \arg \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} - \frac{1}{\beta} r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y}) \right] \\ &= \arg \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp(r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})/\beta)} - \log Z(\mathbf{x}) \right] \\ &= \arg \min_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp(r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})/\beta)} \right],\end{aligned}$$

where  $\frac{1}{Z(\mathbf{x})} = \sum_{\mathbf{y}} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp \frac{r_{\phi^{\dagger}}(\mathbf{x}, \mathbf{y})}{\beta}$ .

And therefore, the optimal value is 0 and optimal solution is

$$\pi^*(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \exp \frac{r_{\phi^*}(\mathbf{x}, \mathbf{y})}{\beta}.$$

Now we can express the unknown score function  $r()$  in terms of optimal solution  $\pi^*$ , hence allow us to reduce the unknown to only  $\pi^*$ .

$$r_{\phi^*}(\mathbf{x}, \mathbf{y}) = \beta \log \frac{\pi^*(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} + \beta \log Z(\mathbf{x})$$



Then with the preference model, we can derive the MLE objective to find that optimal  $\pi^*$ .

- ▶ Under the Bradley-Terry model, observing dataset  $[(\mathbf{x}_i, \mathbf{y}_{1i}, \mathbf{y}_{2i})]_1^n$ , the MLE objective is

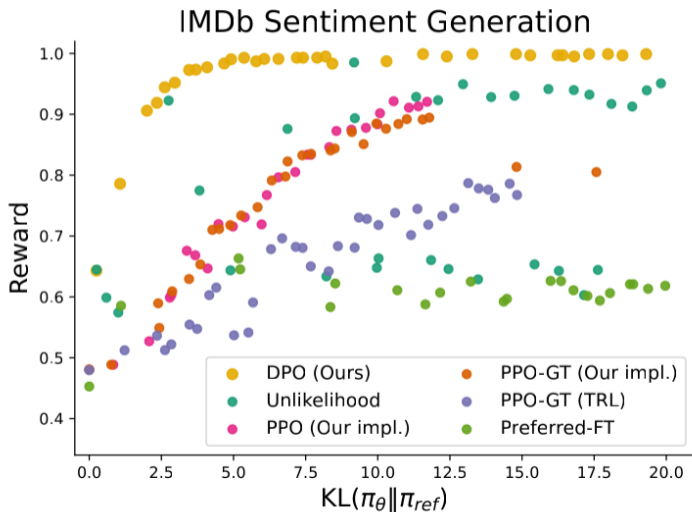
$$\begin{aligned} & \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mathcal{D}} \left[ I[\mathbf{y}_1 \succ \mathbf{y}_2] \sigma(r_\phi(\mathbf{x}, \mathbf{y}_2) - r_\phi(\mathbf{x}, \mathbf{y}_1)) + I[\mathbf{y}_2 \succ \mathbf{y}_1] \sigma(r_\phi(\mathbf{x}, \mathbf{y}_2) - r_\phi(\mathbf{x}, \mathbf{y}_1)) \right] \\ &= \mathbb{E}_{\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2 \sim \mathcal{D}} \left[ I[\mathbf{y}_1 \succ \mathbf{y}_2] \sigma \left( \beta \log \frac{\pi_\phi(\mathbf{y}_1 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1 | \mathbf{x})} - \beta \log \frac{\pi_\phi(\mathbf{y}_2 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2 | \mathbf{x})} \right) \right. \\ & \quad \left. + I[\mathbf{y}_2 \succ \mathbf{y}_1] \sigma \left( \beta \log \frac{\pi_\phi(\mathbf{y}_2 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_2 | \mathbf{x})} - \beta \log \frac{\pi_\phi(\mathbf{y}_1 | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_1 | \mathbf{x})} \right) \right]. \end{aligned}$$

In other words, we are parameterizing the unknown score function

$r(\mathbf{x}, \mathbf{y}) = \log \pi_\theta(\mathbf{x}, \mathbf{y}) - \log \pi_{\text{ref}}(\mathbf{x}, \mathbf{y})$  to guarantee that the optimal solution of problem (1) is  $\pi_\theta$ .

## Control setting

We want to finetune a LM model such that it always produce positive reviews.



## Control setting - My try

- ▶ Dataset: IMDB,  $\sim 20k$  reviews
- ▶ True score function is given by a sentiment classifier (a pretrained large network)
- ▶  $\pi_{\text{ref}}$ : Fine-tuning gpt2-large (1.4B params) on unlabeled IMDB
- ▶ For PPO, we provide the true score function.
- ▶ For DPO, given a prompt, we sample 4 responses for each prompt, and create 6 preference pairs.

Table: About an hour training for each method

	$\pi_{\text{ref}}$	$\pi_{\text{ppo}}$	$\pi_{\text{dpo}}$
Sentiment score	0.625	0.86	0.99
KL	0.	1.7	-26.6

## Result on other tasks

# Extensions

- ▶ Assuming preference pairs are noisy due to annotator's imperfection,

$$z \sim \text{Bern}(\sigma(r(\mathbf{x}, \mathbf{y}_1) - r(\mathbf{x}, \mathbf{y}_2)))$$

$$\ell \sim \text{Pr}(\ell' | z)$$

- ▶ In [Christiano et al. 2017], some pairs annotations are just uniformly selected  $\Rightarrow$  outliers.
- ▶ Instead of pairwise preferences, we can consider a best-choice preferences: Given a prompt  $\mathbf{x}$  and  $L$  responses, the label is the best response. [Ziegler et al. 2019].
- ▶ Assuming existence of score function might not hold in general
- ▶ What about  $D_{\text{kl}}(\pi_{\text{ref}} \parallel \pi_{\theta})$

# Preference Optimization with the Pairwise Cringe Loss

Jing Xu et al. “Some things are more cringe than others: Preference optimization with the pairwise cringe loss”. In: *arXiv preprint arXiv:2312.16682* [2023] Alignment samples can be in different forms:

- ▶ Supervised setting:  $(\mathbf{x}, \mathbf{y})$
- ▶ Binary feedback:  $(\mathbf{x}^+, \mathbf{y}^+, \mathbf{x}^-, \mathbf{y}^-)$
- ▶ Binary preference:  $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$

Cringe loss is originally applied to Binary feedback data:

$$\mathcal{L}_{\text{BIN}}(\mathbf{x}^-, \mathbf{y}^-, \mathbf{x}^+, \mathbf{y}^+) = \mathcal{L}_{\text{CE}} + \mathcal{L}_{\text{Cr}}$$

$$\mathcal{L}_{\text{CE}}(\mathbf{x}^+, \mathbf{y}^+) = -\log \Pr(\mathbf{y}^+ | \mathbf{x}^+)$$

$$\mathcal{L}_{\text{Cr}}(\mathbf{x}^-, \mathbf{y}^-) = -\log \sum_t \log \frac{\exp(s_t^*)}{\exp(s_t^*) + \exp(s_t[y_t^-])},$$

where we feed the prompt  $\mathbf{x}^-$  to the model, and ask it to generate an output of length  $T$ :

- ▶ At the  $t$ -th token, we select top  $k$  tokens according model's prob output  $s_t^1, \dots, s_t^k$ .
- ▶ Normalizing probability over these tokens by applying softmax function.
- ▶ Sample an index  $z \sim \text{Categorical}(s_t^1, \dots, s_t^k)$ ,  $z \in [k]$ .
- ▶  $s_t^* = s_t^z$ .

## Apply Cringe Loss to Pairwise Preference data

They propose to use the following loss on pairwise preference data

$$\mathcal{L}_{\text{Pair}}(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = g(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) \mathcal{L}_{\text{BIN}}(\mathbf{x}, \mathbf{y}_1, \mathbf{x}, \mathbf{y}_2),$$

where

$$g(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = \sigma(b - M(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)),$$

$$M(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) = \log \Pr(\mathbf{y}_1 | \mathbf{x}) - \log \Pr(\mathbf{y}_2 | \mathbf{x}).$$



# Result

Table 1: AlpacaFarm evaluation results (LLM evaluation), using human preference data and reward model (where applicable) for training. (\*=average of 3 seeds). <sup>1</sup>PPO with human preferences was trained by [Dubois et al. \(2023\)](#); we just evaluated the model.

METHOD	WIN RATE (%)
<i>Results reported by <a href="#">Dubois et al. (2023)</a></i>	
LLAMA 7B	11.3
SFT 10k	36.7
SFT 52k	39.2
<i>Experiments reported in this paper:</i>	
BINARY CRINGE	47.7*
PPO <sup>1</sup>	48.5*
DPO	50.2*
PAIRWISE CRINGE	54.7*

Figure: Image