

The Crucial Role of Samplers in Online Direct Preference Optimization

Ruizhe Shi*, Runlong zhou*, Simon S. Du

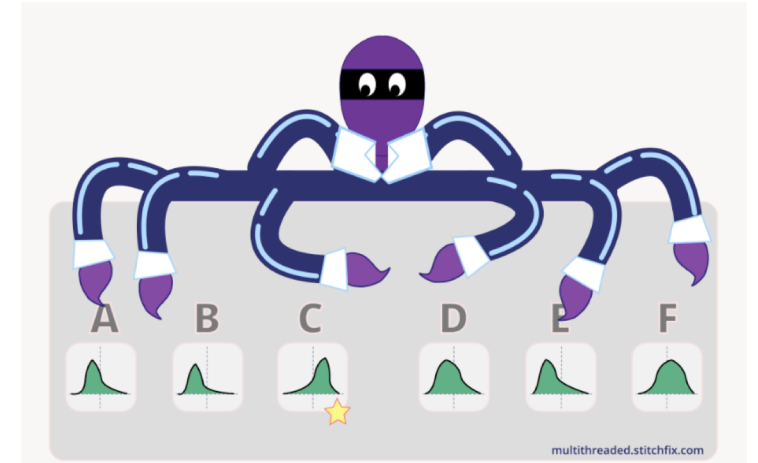
Mar 9 2025

ICLR 2025

Bandit view of language model alignment

- **Prompt (State):** user queries (x)
- **Response (Action):** language model generation (y)
- **Reward function:** $r(x, y) \in [0,1]$

(in this project we only study MAB)



Picture from
<https://multithreaded.stitchfix.com/blog/2020/08/05/bandits/>

Policy

- A **tabular softmax** policy π_θ for MABs satisfies

$$\pi_\theta(y) = \frac{e^{\theta_y}}{\sum_{y'} e^{\theta_{y'}}$$

Preference-based RL

- A **preference** model $p^*(y_1 \succ y_2)$ indicating the probability that y_1 is preferred over y_2
- After choosing a **pair** of arms (y_1, y_2) , observe a sample $p \sim \text{Bernoulli}(p^*(y_1 \succ y_2))$
- BT preference model

$$p^*(y_1 \succ y_2) = \sigma(r(y_1) - r(y_2)) = \frac{e^{r(y_1)}}{e^{r(y_1)} + e^{r(y_2)}}$$

Motivation: how fast can data help DPO converge

- Human preference dataset $\mathcal{D} = \left\{ \left(y_w^{(i)}, y_l^{(i)} \right) \right\}_{i=1}^N$
 - In the i th sample, $y_w^{(i)}$ is preferred over $y_l^{(i)}$
- DPO (a popular alignment algorithm)
 - $\mathcal{L}_\pi(\theta) = -\frac{1}{N} \sum_{i=1}^N \log \sigma \left(\beta \log \frac{\pi_\theta(y_w^{(i)})}{\pi_{\text{ref}}(y_w^{(i)})} - \beta \log \frac{\pi_\theta(y_l^{(i)})}{\pi_{\text{ref}}(y_l^{(i)})} \right)$
 - Closed-form solution: $\pi^*(y) = \frac{1}{Z} \pi_{\text{ref}}(y) e^{r(y)/\beta}$
- We want to study: *how fast can DPO converge to optimality with different sampling distributions on data?*

Motivation: how fast can data help DPO converge

How fast can $r(y) - r(y') - \beta \log \frac{\pi_{\theta^{(t)}}(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_{\theta^{(t)}}(y')}$ *converge to 0, for* $\forall y, y' \in \mathcal{Y}$?

$\underbrace{\hspace{15em}}_{=: \delta(y, y'; \theta^{(t)})}$

Ideal Case: Exact DPO

- Suppose we have two **sampling policies** π^{s1} for y_1 and π^{s2} for y_2
- Define sampling probability

Stop gradient

$$\pi^s(y, y') := \text{sg} (\pi^{s1}(y)\pi^{s2}(y') + \pi^{s1}(y')\pi^{s2}(y))$$

- Exact DPO loss function

$$\mathcal{L}_{\text{DPO}}(\theta) := - \sum_{y, y' \in \mathcal{Y}} \pi^s(y, y') p^*(y > y') \log \sigma \left(\beta \log \frac{\pi_\theta(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_\theta(y')} \right)$$

- Policy update

$$\theta^{(t+1)} = \theta^{(t)} - \eta \alpha(\pi^{s1}, \pi^{s2}) \nabla_\theta \mathcal{L}_{\text{DPO}}(\theta^{(t)})$$

Sampling coefficient determined by samplers

Ideal Case: Exact DPO

- Mixture of samplers

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} \left(\alpha_1 \mathcal{L}_1(\theta^{(t)}) + \alpha_2 \mathcal{L}_2(\theta^{(t)}) \right)$$

- Central to our design

Practical Case: Empirical DPO

- No access to exact gradients

$$\theta^{(t+1)} = \theta^{(t)} - \eta G^{(t)}$$

where $G_y^{(t)}$ is a random variable that

$$\frac{1}{\beta A} \left(G_y^{(t)} - \alpha(\pi^{s1}, \pi^{s2}) \nabla_{\theta_y} \mathcal{L}(\theta^{(t)}) \right) \sim \text{sub-Gaussian}(\sigma^2)$$

- Mixture of samplers

$$\frac{1}{\beta A} \left(G_y^{(t)} - \nabla_{\theta_y} \left(\alpha_1 \mathcal{L}_1(\theta^{(t)}) + \alpha_2 \mathcal{L}_2(\theta^{(t)}) \right) \right) \sim \text{sub-Gaussian}(\sigma^2)$$

Main results

- Uniform sampler (*vanilla*) $\pi^{s1}(\cdot) = \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y})$

$$\left| \delta(y, y'; \theta^{(T)}) \right| \leq 0.588^T, \quad \forall y, y' \in \mathcal{Y}$$

$$\max_{y, y' \in \mathcal{Y}} \left| \delta(y, y'; \theta^{(T)}) \right| \geq \gamma^T \quad \text{linear convergence!}$$

determined by initialization

- Policy-difference guided sampler (*ours*)

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot) = \text{Uniform}(\mathcal{Y}), \\ \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}), \end{cases} \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi(\cdot)/\pi_{\text{ref}}(\cdot))^\beta \\ \pi^{s2}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi_{\text{ref}}(\cdot)/\pi(\cdot))^\beta \end{cases}$$

$$\left| \delta(y, y'; \theta^{(T)}) \right| \leq 0.611^{2^T - 1}, \quad \forall y, y' \in \mathcal{Y} \quad \text{quadratic convergence!}$$

Regime 1: Known Reward

Not practical, only for proof of idea

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot) = \text{Uniform}(\mathcal{Y}) \\ \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}) \end{cases}, \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot \exp(r(\cdot)) \\ \pi^{s2}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot \exp(-r(\cdot)) \end{cases},$$

- Sampling coefficient $\alpha_1 = |\mathcal{Y}|^2$, $\alpha_2 = \sum_{y, y'} \exp(r(y) - r(y'))$
- Upper bound

Quadratic convergence!

$$\left| \delta(y, y'; \theta^{(T)}) \right| \leq 0.5^{2^T - 1}, \quad \forall y, y' \in \mathcal{Y}$$

Intuition

$$\propto \delta(y, y''; \theta^{(t)}) - \delta(y', y''; \theta^{(t)}) + \mathcal{O}(\delta^2) = \delta(y, y'; \theta^{(t)}) + \mathcal{O}(\delta^2)$$

$$\delta(y, y'; \theta) := r(y) - r(y') - \beta \log \frac{\pi_\theta(y) \pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y) \pi_\theta(y')}$$

We care about its convergence

- Recall update

$$\delta(y, y'; \theta^{(t+1)}) = \delta(y, y'; \theta^{(t)}) - \eta\beta \sum_{y'' \in \mathcal{Y}} \left[\pi^s(y, y'') \Delta(y, y''; \theta^{(t)}) - \pi^s(y', y'') \Delta(y', y''; \theta^{(t)}) \right]$$

where $\Delta(y, y'; \theta) := \sigma(r(y) - r(y')) - \sigma \left(\beta \log \frac{\pi_\theta(y) \pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y) \pi_\theta(y')} \right)$

- Taylor expansion at $r(y_1) - r(y_2)$ and setting $\pi^s(y_1, y_2) \propto 1/\sigma'(r(y_1) - r(y_2))$ gives

$$\pi^s(y, y'') \Delta(y, y''; \theta^{(t)}) - \pi^s(y', y'') \Delta(y', y''; \theta^{(t)}) = \text{constant} \cdot \delta(y, y'; \theta^{(t)}) + \text{quadratic term}$$

Regime 1: Known Reward

- The choice of η eliminates the linear term:

$$\begin{aligned} \delta(a, a'; \theta^{(t+1)}) &= (1 - \eta\beta^2 A)\delta(a, a'; \theta^{(t)}) \\ &+ \frac{\eta\beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_R(a, a''; \theta^{(t)}))}{\sigma'(r(a) - r(a''))} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_R(a', a''; \theta^{(t)}))}{\sigma'(r(a') - r(a''))} \delta(a', a''; \theta^{(t)})^2 \right) \end{aligned}$$

- Bounding $\sigma'' \leq \frac{1}{6\sqrt{3}} < 0.097$ and $\sigma' \geq \sigma'(1) > 0.196$ gives
$$|\delta(y, y'; \theta^{(t+1)})| < 0.5 \max_{a, a'} \delta(a, a'; \theta^{(t)})^2$$

Regime 2: Online Sampler

Current policy

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot) = \text{Uniform}(\mathcal{Y}), \\ \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}), \end{cases} \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi(\cdot)/\pi_{\text{ref}}(\cdot))^\beta \\ \pi^{s2}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi_{\text{ref}}(\cdot)/\pi(\cdot))^\beta \end{cases}$$

- $\textcircled{2}$ equivalent to $\pi^{s1} \propto \exp(\beta(\theta - \theta_{\text{ref}}))$, $\pi^{s2} \propto \exp(\beta(\theta_{\text{ref}} - \theta))$
- Sampling coefficient $\alpha_1 = |\mathcal{Y}|^2$, $\alpha_2 = \sum_{y, y'} \left(\frac{\pi(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi(y')} \right)^\beta$
- Upper bound

Quadratic convergence!

$$\left| \delta(y, y'; \theta^{(T)}) \right| \leq 0.611^{2^T - 1}, \quad \forall y, y' \in \mathcal{Y}$$

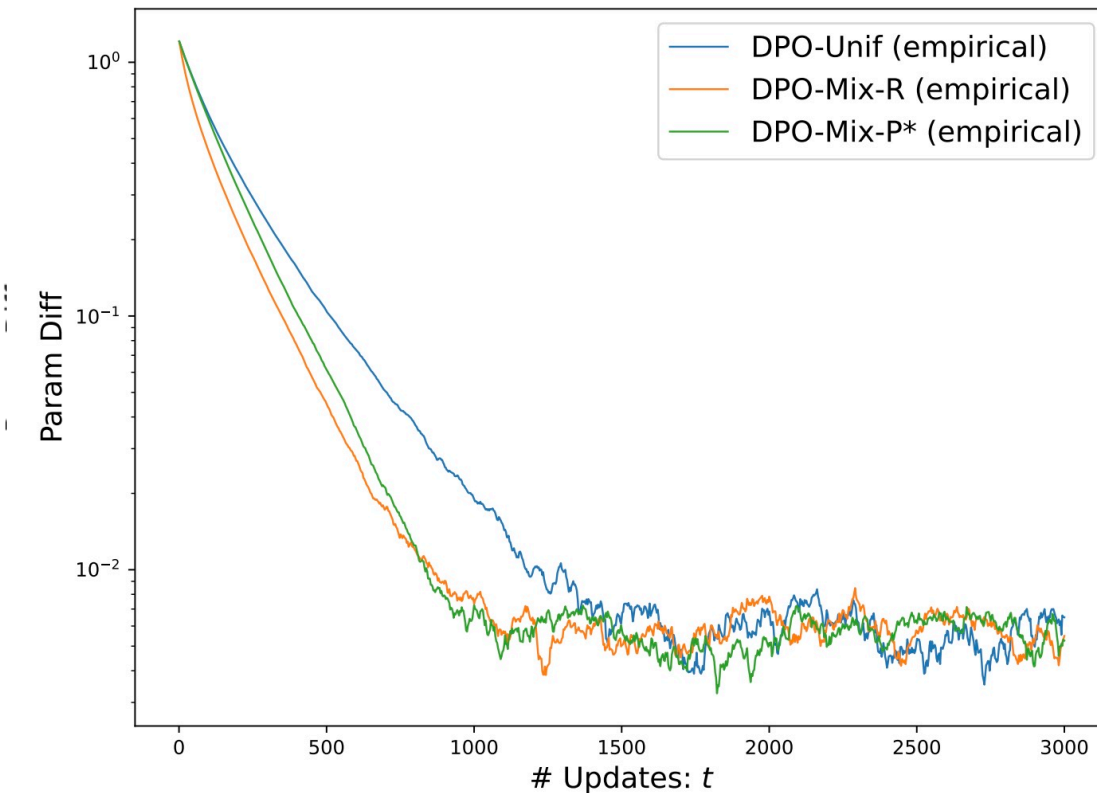
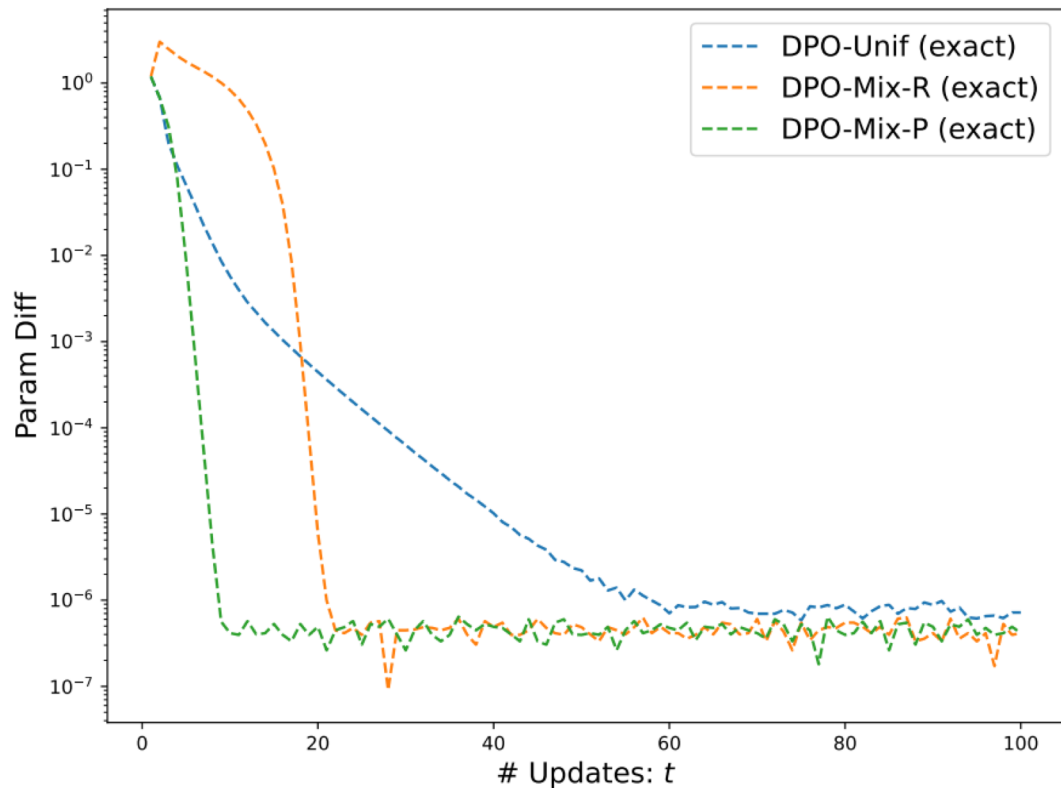
Regime 2: Online Sampler

- Taylor expansion at $\beta \log \frac{\pi(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi(y')}$

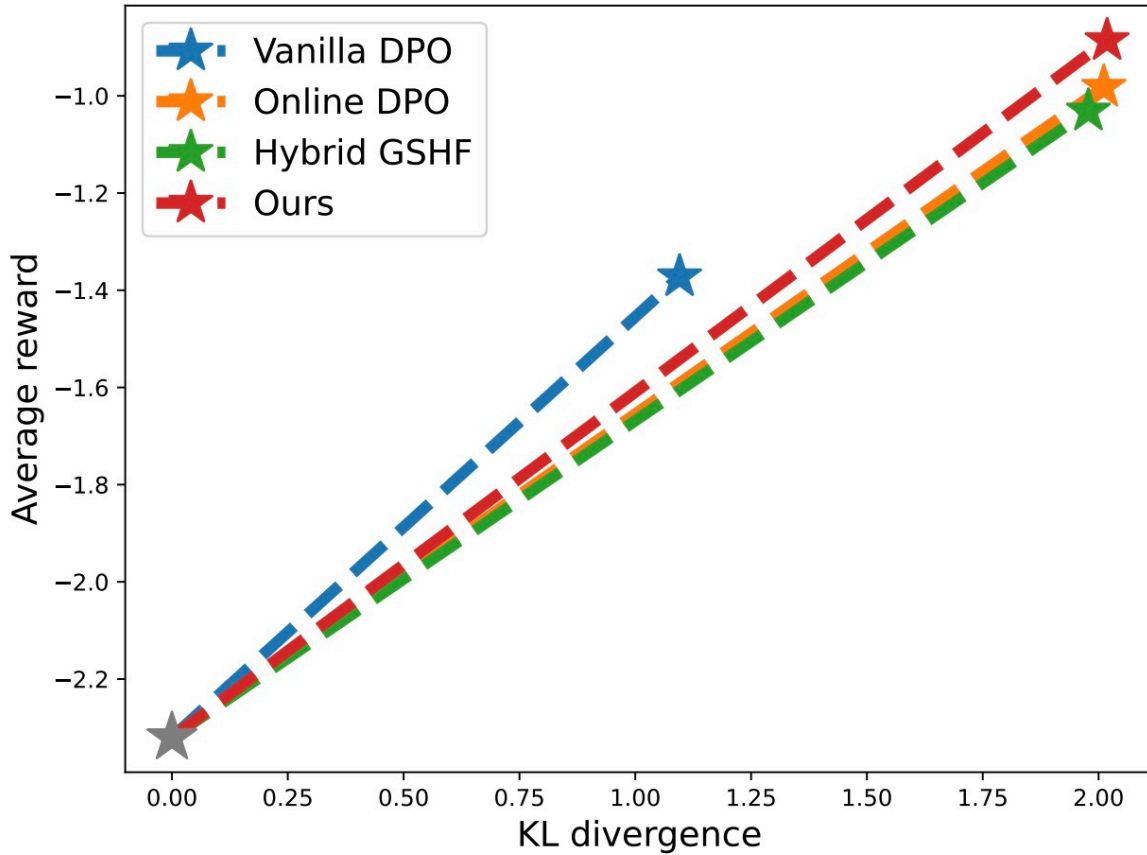
$$\delta(a, a'; \theta^{(t+1)}) = (1 - \eta\beta^2 A)\delta(a, a'; \theta^{(t)}) - \frac{\eta\beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_P(a, a''; \theta^{(t)}))}{\sigma'(\beta(\theta_a - \theta_{a''})^{(t)})} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_P(a', a''; \theta^{(t)}))}{\sigma'(\beta(\theta_{a'} - \theta_{a''})^{(t)})} \delta(a', a''; \theta^{(t)})^2 \right)$$

- Then same as regime 1.

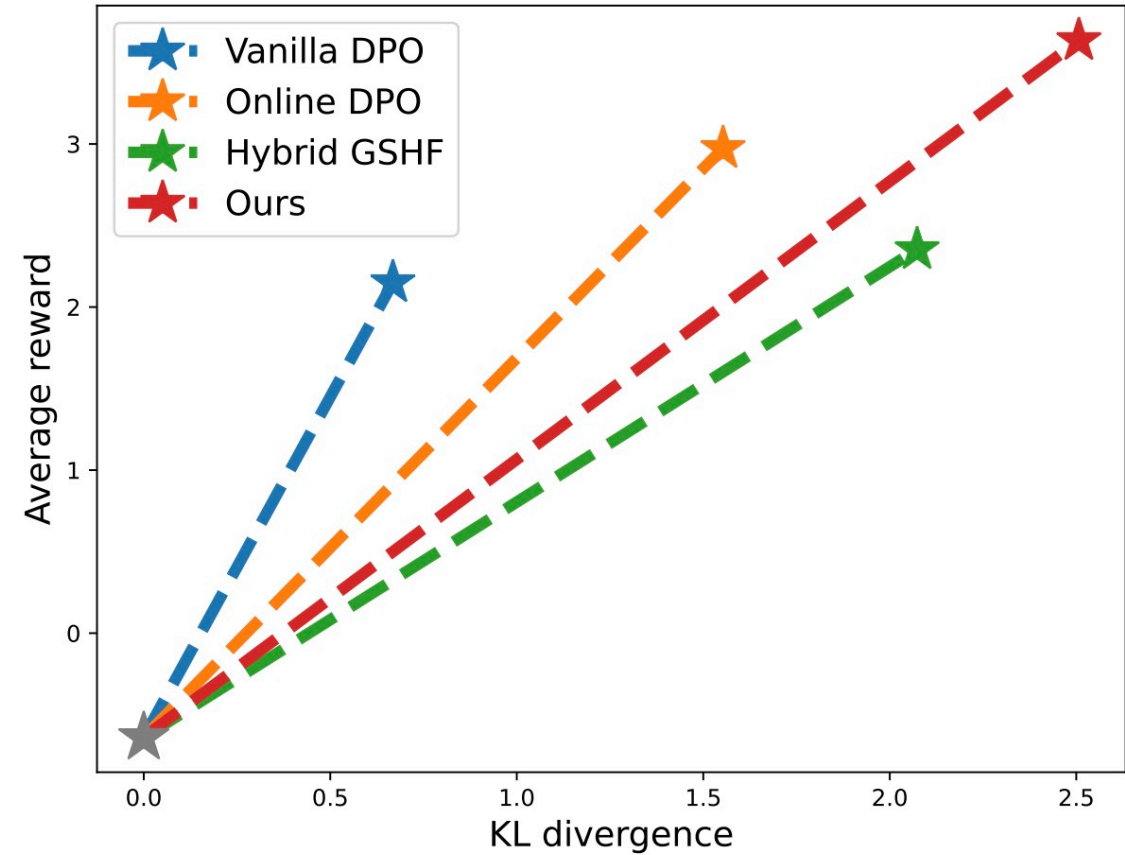
Numerical Simulations



LM Implementation



Safe-RLHF



Iterative Prompt

Thank You