

# SQUAREDPO: Displacement-Resistant Extensions of DPO with Nonconvex $f$ -Divergences

Idan Pipano   Shoham Sabach   Kavosh Asadi  
Mohammad Ghavamzadeh







Alignment: we start with an LLM  $\pi_{\text{ref}}$ , and wish to align it using a dataset of triplets  $(x, y_w, y_l)$ .

- ▶ In RLHF, there are two steps:

- ▶ In RLHF, there are two steps:
  - ▶ we first learn a reward model  $r_\phi$

$$\min_{\phi} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [-\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$$

- ▶ In RLHF, there are two steps:

- ▶ we first learn a reward model  $r_\phi$

$$\min_{\phi} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [-\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$$

- ▶ and then train the LLM  $\pi_\theta$

$$\max_{\theta} \mathbb{E}_x [\mathbb{E}_{y \sim \pi_\theta(\cdot | x)} [r_\phi(x, y)] - \beta D_{\text{KL}}[\pi_\theta(\cdot | x) \parallel \pi_{\text{ref}}(\cdot | x)]]$$

- ▶ In RLHF, there are two steps:

- ▶ we first learn a reward model  $r_\phi$

$$\min_{\phi} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [-\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$$

- ▶ and then train the LLM  $\pi_\theta$

$$\max_{\theta} \mathbb{E}_x [\mathbb{E}_{y \sim \pi_\theta(\cdot | x)} [r_\phi(x, y)] - \beta D_{\text{KL}}[\pi_\theta(\cdot | x) \parallel \pi_{\text{ref}}(\cdot | x)]]$$

- ▶ DPO collapses the two-step RLHF process into a single loss
  - ▶ by substituting the optimal solution of the second RLHF loss into the first one.

- ▶ A recent generalization of DPO is  $f$ -DPO, obtained by replacing  $D_{\text{KL}}$  with the more general  $D_f$ .

- ▶ A recent generalization of DPO is  $f$ -DPO, obtained by replacing  $D_{\text{KL}}$  with the more general  $D_f$ .
- ▶ **Our first contribution:**

## Theorem 1

$f$  can be used in  $f$ -DPO if and only if  $\lim_{t \rightarrow 0^+} f'(t) = -\infty$

- ▶ A recent generalization of DPO is  $f$ -DPO, obtained by replacing  $D_{\text{KL}}$  with the more general  $D_f$ .
- ▶ **Our first contribution:**

### Theorem 1

$f$  can be used in  $f$ -DPO if and only if  $\lim_{t \rightarrow 0^+} f'(t) = -\infty$

- ▶ In particular,  $f$  does not have to be convex!

- ▶ Training an LLM  $\pi_\theta$  with DPO causes both  $\pi_\theta(y_l | x)$  and  $\pi_\theta(y_w | x)$  (!) to decrease.

- ▶ Training an LLM  $\pi_\theta$  with DPO causes both  $\pi_\theta(y_l | x)$  and  $\pi_\theta(y_w | x)$  (!) to decrease.
  - ▶ We call this phenomenon **likelihood displacement**.

- ▶ Training an LLM  $\pi_\theta$  with DPO causes both  $\pi_\theta(y_l | x)$  and  $\pi_\theta(y_w | x)$  (!) to decrease.
  - ▶ We call this phenomenon **likelihood displacement**.
- ▶ Widely observed empirically, poorly understood theoretically.

- ▶ Training an LLM  $\pi_\theta$  with DPO causes both  $\pi_\theta(y_l | x)$  and  $\pi_\theta(y_w | x)$  (!) to decrease.
  - ▶ We call this phenomenon **likelihood displacement**.
- ▶ Widely observed empirically, poorly understood theoretically.
- ▶ **Our second contribution:**

## Theorem 2

A necessary condition on  $f$  to mitigate likelihood displacement in  $f$ -DPO is  $\arg \min_t f(t) \geq 1$ .

- ▶ Training an LLM  $\pi_\theta$  with DPO causes both  $\pi_\theta(y_l | x)$  and  $\pi_\theta(y_w | x)$  (!) to decrease.
  - ▶ We call this phenomenon **likelihood displacement**.
- ▶ Widely observed empirically, poorly understood theoretically.
- ▶ **Our second contribution:**

## Theorem 2

A necessary condition on  $f$  to mitigate likelihood displacement in  $f$ -DPO is  $\arg \min_t f(t) \geq 1$ .

- ▶ In particular, the  $f$  of DPO satisfies  $\arg \min_t f(t) = e^{-1}$ ; hence, one can expect likelihood displacement in DPO.

- ▶ Our theory imposes two desiderata on  $f$  for  $f$ -DPO.
  - ▶  $\lim_{t \rightarrow 0^+} f'(t) = -\infty$
  - ▶  $\arg \min_t f(t) \geq 1$

- ▶ Our theory imposes two desiderata on  $f$  for  $f$ -DPO.
  - ▶  $\lim_{t \rightarrow 0^+} f'(t) = -\infty$
  - ▶  $\arg \min_t f(t) \geq 1$
- ▶  $f(t) = (\ln t)^2$  is an example of a function that satisfies them.

- ▶ Our theory imposes two desiderata on  $f$  for  $f$ -DPO.
  - ▶  $\lim_{t \rightarrow 0^+} f'(t) = -\infty$
  - ▶  $\arg \min_t f(t) \geq 1$
- ▶  $f(t) = (\ln t)^2$  is an example of a function that satisfies them.
- ▶ We call the  $f$ -DPO loss with this  $f$  SQUAREDPO.

- Our method: SQUAREDPO =  $f$ -DPO with  $f(t) = (\ln t)^2$

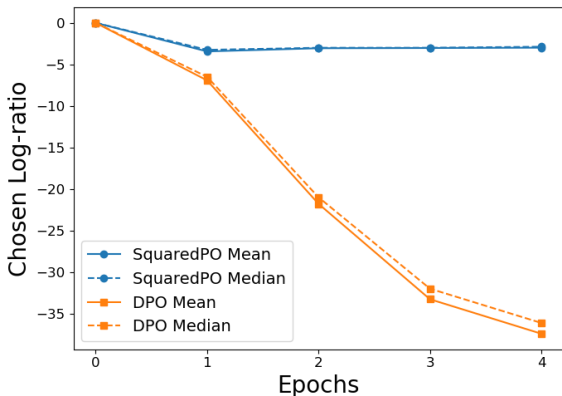


Figure: Chosen Log-ratio is  $\log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)}$ .

- Our method: SQUAREDPO =  $f$ -DPO with  $f(t) = (\ln t)^2$

Table: Win-rate on the validation set against  $\pi_{\text{ref}}$ .

| Epochs | SQUAREDPO (%)  | DPO (%)        |
|--------|----------------|----------------|
| 1      | $50.8 \pm 0.7$ | $51.8 \pm 1.0$ |
| 2      | $50.6 \pm 1.1$ | $45.0 \pm 1.1$ |
| 4      | $51.0 \pm 0.7$ | $34.7 \pm 1.3$ |

- Our method: SQUAREDPO =  $f$ -DPO with  $f(t) = (\ln t)^2$

Table: Win-rate on the validation set against  $\pi_{\text{ref}}$ .

| Epochs | SQUAREDPO (%)  | DPO (%)        |
|--------|----------------|----------------|
| 1      | $50.8 \pm 0.7$ | $51.8 \pm 1.0$ |
| 2      | $50.6 \pm 1.1$ | $45.0 \pm 1.1$ |
| 4      | $51.0 \pm 0.7$ | $34.7 \pm 1.3$ |

- SQUAREDPO achieves performance comparable to DPO while being more robust to the number of training epochs.

# Thank You for Listening!



The Paper 