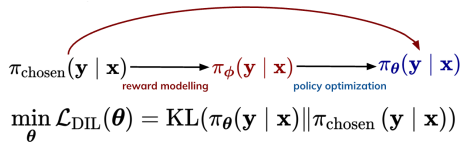


Insight: Directly conducting imitation learning

Without Bradley-Terry assumption



Framework: Unlock family of Alignment Algorithms

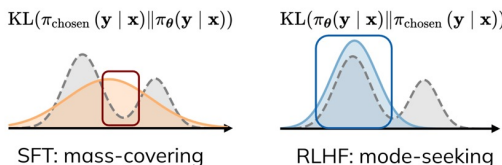
$$\max_{\theta} \mathbb{E}_{\pi_{\theta}(\mathbf{y}|\mathbf{x})} \left[\log \frac{\pi_{\text{chosen}}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} - \log \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right] = \mathbb{E}_{\pi_{\theta}(\mathbf{y}|\mathbf{x})} [\log r(\mathbf{x}, \mathbf{y})] - \text{KL}(\pi_{\theta}(\mathbf{y} | \mathbf{x}) || \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))$$

The reward is about the density ratio

Table 1: Summary of the variants of DIL with different h -functions for Bregman divergence: $\mathcal{L}_{\text{DIL}}(\theta) = \mathbb{E}_{\pi_{\text{chosen}}(\mathbf{y}|\mathbf{x})} [f_1(f_{\theta})] + \mathbb{E}_{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})} [f_2(f_{\theta})]$ as a function of log ratio $f_{\theta} = \log(\pi_{\theta}(\mathbf{y} | \mathbf{x}) / \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}))$.

h -Bregman Density Ratio Estimation	h -function	$f_1(f_{\theta})$	$f_2(f_{\theta})$
LSIF (Kumar et al., 2009)	$h(r) = (r - 1)^2/2$	$-e^{f_{\theta}}$	$\frac{1}{2}e^{2f_{\theta}}$
BCE (Hartle et al., 2009)	$h(r) = r \log r - (r+1) \log(r+1)$	$\log(1 + e^{-f_{\theta}})$	$\log(1 + e^{f_{\theta}})$
UKL (Nguyen et al., 2010)	$h(r) = r \log r - r$	$-f_{\theta}$	$e^{f_{\theta}}$

SFT vs RLHF: Different Directions



On a Connection Between Imitation Learning and RLHF

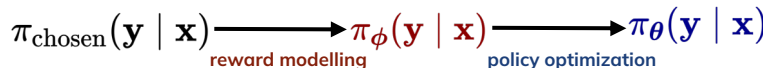
RLHF is Barely RL, and Secretly Performs Imitation Learning

Key Finding: RLHF performs imitation learning, not RL!

Proposition 1. Suppose the chosen response distribution $\pi_{\text{chosen}}(\mathbf{y} | \mathbf{x})$, the EBM $\pi_{\phi}(\mathbf{y} | \mathbf{x})$, and the model $\pi_{\theta}(\mathbf{y} | \mathbf{x})$. KL-regularized RLHF with $\beta = 1$ can be viewed as the following problem:

$$\min_{\pi_{\theta}} \text{KL}(\pi_{\theta} || \pi_{\phi}^*) \quad \text{s.t.} \quad \pi_{\phi}^* = \arg \min_{\pi_{\phi}} \text{KL}(\pi_{\text{chosen}} || \pi_{\phi}),$$

where $\pi_{\text{chosen}}(\mathbf{y} | \mathbf{x}) = \pi_{\phi}(\mathbf{y} | \mathbf{x}) = \pi_{\theta}(\mathbf{y} | \mathbf{x})$ is the equilibrium.



Result: DIL Outperforms SOTAs across Benchmarks

Model (↓) / Benchmark (→)	MLLM-PRO	BBH	MUSR	MATH	GSM8K	ARC
SFT	31.00	46.16	41.27	3.70	46.32	60.15
DPO (Rafailov et al., 2024)	31.58	47.80	40.48	4.53	38.67	64.42
SLiC (Zhao et al., 2023)	31.11	46.53	40.55	3.92	48.82	61.43
f-DPO (Wang et al., 2024a)	30.85	47.55	40.39	4.37	39.55	62.85
IPO (Azar et al., 2024)	30.18	46.78	39.58	4.02	22.67	62.88
KTO (Ehlayaraj et al., 2024)	31.16	47.92	40.24	4.13	38.99	63.17
CPO (Xu et al., 2024b)	30.95	47.17	41.59	4.25	46.93	61.69
SimPO (Meng et al., 2024)	31.61	48.38	40.08	4.23	31.54	65.19
DIL w/ LSIF	32.22	48.78	42.75	4.68	48.98	65.37

Dataset (→)	TL;DR Summarization			Anthropic-HH		
Method (↓) / Metric (→)	vs SFT	vs Chosen	Average	vs SFT	vs Chosen	Average
DPO (Rafailov et al., 2024)	71.22	57.58	64.40	69.32	59.35	64.34
SLiC (Zhao et al., 2023)	68.61	55.72	62.17	65.52	57.71	61.62
f-DPO (Wang et al., 2024a)	66.19	51.37	58.78	60.21	52.38	56.30
IPO (Azar et al., 2024)	72.17	56.51	64.34	63.19	55.12	59.16
CPO (Xu et al., 2024b)	73.13	58.89	66.01	72.30	63.39	67.86
SimPO (Meng et al., 2024)	69.71	54.38	62.05	67.85	57.51	62.68
DIL w/ LSIF	75.47	60.25	67.86	73.32	65.02	69.17

