

Policy-labeled Preference Learning: Is Preference Enough for RLHF?

Taehyun Cho^{1,†} Seokhun Ju^{1,†} Seungyub Han¹
Dohyeong Kim¹ Kyungjae Lee² Jungwoo Lee¹

¹Seoul National University, ²Korea University

[†]Equal Contribution

ICML (spotlight)
June 20, 2025



Seoul National University
Cognitive Machine Learning Lab.

KOREA
UNIVERSITY

- Reinforcement Learning from Human Feedback (RLHF) enables agents to align with human goals using human preferences.

- Reinforcement Learning from Human Feedback (RLHF) enables agents to align with human goals using human preferences.
- Existing RLHF methods often assume as if trajectories are generated by optimal policies π^* .

- Reinforcement Learning from Human Feedback (RLHF) enables agents to align with human goals using human preferences.
- Existing RLHF methods often assume as if trajectories are generated by optimal policies π^* .
- This leads to **likelihood mismatch** in offline settings due to environmental stochasticity and diverse behavior policies.

- Reinforcement Learning from Human Feedback (RLHF) enables agents to align with human goals using human preferences.
- Existing RLHF methods often assume as if trajectories are generated by optimal policies π^* .
- This leads to **likelihood mismatch** in offline settings due to environmental stochasticity and diverse behavior policies.
- Direct Preference Optimization (DPO) removes the need for explicit rewards, but fails to address this mismatch.

- 1 We propose **Policy-labeled Preference Learning (PPL)**, a regret-based framework for RLHF, which explicitly models the **behavior policy** associated with preference data.
- 2 We introduce **contrastive KL regularization** to correct for likelihood mismatch.
- 3 PPL shows superior performance on MetaWorld offline tasks and is competitive in online RLHF.

Score-based Preference Model

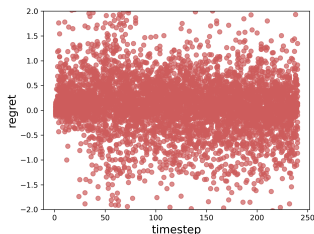
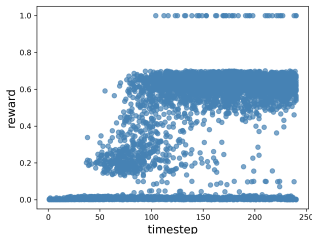
Table 1: Comparison for different preference models under PbRL framework.

Algorithm	Score Function	Direct Preference Optimization	Likelihood Matching
PEBBLE (Lee et al., 2021)	$r_\psi(s_t, a_t)$	✗	✗
DPO (Rafailov et al., 2024b)	$\log \pi_\psi(y s) / \pi_{\text{ref}}(y s)$	✓	✗
DPPO (An et al., 2023)	$-\mathbb{E}_{a \sim \pi_\psi(\cdot s_t)} [\ a - a_t\ _2]$	✓	✗
CPL (Hejna et al., 2023)	$Q^{\pi_\psi}(s_t, a_t) - V^{\pi_\psi}(s_t)$	✓	✗
PPL [Ours]	$-(V^{\pi_\psi}(s_t) - Q^\pi(s_t, a_t))$	✓	✓

Model prediction: $P_{S_\psi}[\zeta^+ \succ \zeta^-] = \sigma\left(\sum_{t \geq 0} S_\psi(s_t^+, a_t^+) - S_\psi(s_t^-, a_t^-)\right),$

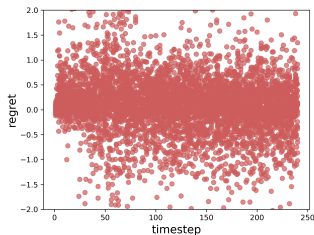
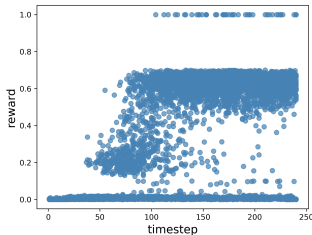
Loss function: $\mathcal{L}(S_\psi; \mathcal{D}) = -\mathbb{E}_{(\zeta^+, \zeta^-) \sim \mathcal{D}} \left[\log P_{S_\psi}[\zeta^+ \succ \zeta^-] \right]$

Reward vs Regret



- **Reward** : Sparse and delayed feedback
- **Negative Regret** : Dense and stepwise feedback

Reward vs Regret



Negative Regret

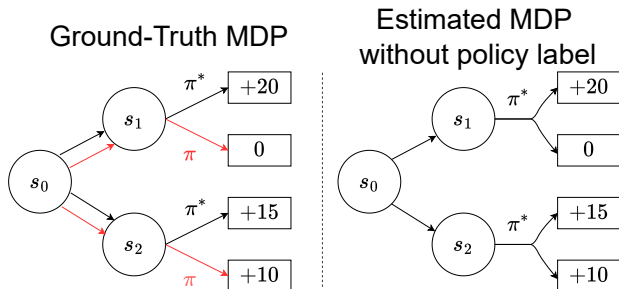
- Performance difference between the **behavior policy** π and the optimal policy π^*
- $-\text{Reg}_{\pi^*}^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi^*}(s)$

From a perspective of regret, existing RLHF/DPO disregards the source of the trajectories, implicitly treating all trajectories as if they were generated by the optimal policy.

From a perspective of regret, existing RLHF/DPO disregards the source of the trajectories, implicitly treating all trajectories as if they were generated by the optimal policy.

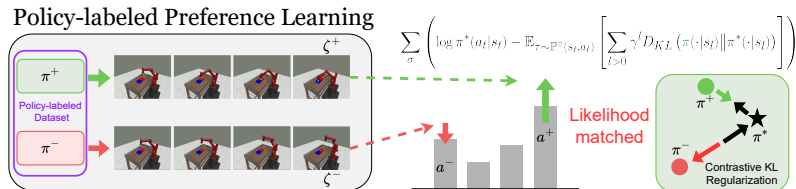
“What impact does this assumption – treating all behavior policies as optimal – have on the regret-based learning process?”

Likelihood Mismatch



- Offline data from π is misinterpreted as from π^* .
- Leads to erroneous preference modeling and degraded performance.

Policy-labeled Preference Learning



- Prior works like CPL use optimal advantage $Q^{\pi^*}(s, a) - V^{\pi^*}(s)$ as preference score.
- But this assumes all data comes from π^* , ignoring suboptimal behavior policies.
- PPL instead uses **negative regret**: $Q^{\pi}(s, a) - V^{\pi^*}(s)$, which incorporates behavior policy.

Sequential Forward KL Divergence

Theorem (Policy Deviation Theorem)

If a policy π^* is α -optimal, then for any policy π ,

$$Q_{*}^{\pi^*}(s, a) - Q_{*}^{\pi}(s, a) = \alpha \bar{D}_{KL}(\pi || \pi^*; s, a)$$

where the **sequential forward KL divergence** is defined as

$$\bar{D}_{KL}(\pi || \pi'; s, a) := \mathbb{E}_{\tau \sim \mathbb{P}_{s,a}^{\pi}} \left[\sum_{l>0} \gamma^l D_{KL}(\pi(\cdot | s_l) || \pi'(\cdot | s_l)) \right].$$

Here, $\mathbb{P}_{s,a}^{\pi}$ is the distribution of trajectories $\tau = (s_0, a_0, \dots, s_l, a_l, \dots)$ generated by policy π and the transition \mathbb{P} , starting at $(s_0, a_0) = (s, a)$.

Sequential Forward KL Divergence

Now we can derive the (negative) regret into policy expression,

$$\begin{aligned} -\text{Reg}_{\pi^*}^{\pi}(s_t, a_t) &:= - \underbrace{V^{\pi^*}(s_t)}_{\text{expected return under } \pi^*} + \underbrace{Q^{\pi}(s_t, a_t)}_{\text{achieved return under } \pi} \\ &= \alpha \left(\underbrace{\log \pi^*(a_t | s_t)}_{\text{increase likelihood}} - \underbrace{\bar{D}_{\text{KL}}(\pi || \pi^*; s_t, a_t)}_{\text{decrease sequential forward KL}} \right). \end{aligned}$$

$$\mathcal{L}_{\text{PPL}}(\pi_{\psi}; \mathcal{D}) = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(- \sum_{t \geq 0} \text{Reg}_{\pi_{\psi}}^{\pi^+}(s_t^+, a_t^+) - \text{Reg}_{\pi_{\psi}}^{\pi^-}(s_t^-, a_t^-) \right) \right]$$

Policy-labeled Preference Learning

$$\mathcal{L}_{\text{PPL}}(\pi_\psi; \mathcal{D}) = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(- \sum_{t \geq 0} \text{Reg}_{\pi_\psi}^{\pi^+}(s_t^+, a_t^+) - \text{Reg}_{\pi_\psi}^{\pi^-}(s_t^-, a_t^-) \right) \right]$$

Substitute $-\text{Reg}_{\pi^*}^{\pi}(s_t, a_t) = \alpha \left(\log \pi^*(a_t | s_t) - \bar{D}_{\text{KL}}(\pi || \pi^*; s_t, a_t) \right)$

Policy-labeled Preference Learning

$$\mathcal{L}_{\text{PPL}}(\pi_\psi; \mathcal{D}) = -\mathbb{E}_{\mathcal{D}} \left[\log \sigma \left(- \sum_{t \geq 0} \text{Reg}_{\pi_\psi}^{\pi^+}(s_t^+, a_t^+) - \text{Reg}_{\pi_\psi}^{\pi^-}(s_t^-, a_t^-) \right) \right]$$

Substitute $-\text{Reg}_{\pi^*}^{\pi}(s_t, a_t) = \alpha \left(\log \pi^*(a_t | s_t) - \bar{D}_{\text{KL}}(\pi || \pi^*; s_t, a_t) \right)$

$$\begin{aligned} & - \sum_{t \geq 0} \text{Reg}_{\pi_\psi}^{\pi^+}(s_t^+, a_t^+) - \text{Reg}_{\pi_\psi}^{\pi^-}(s_t^-, a_t^-) \\ &= \alpha \sum_{t \geq 0} \left(\log \frac{\pi_\psi(a_t^+ | s_t^+)}{\pi_\psi(a_t^- | s_t^-)} \underbrace{- \bar{D}_{\text{KL}}(\pi^+ || \pi_\psi; s_t^+, a_t^+) + \bar{D}_{\text{KL}}(\pi^- || \pi_\psi; s_t^-, a_t^-)}_{\text{contrastive KL regularization } \mathcal{R}(\pi_\psi; \pi^+, \pi^-)} \right) \end{aligned}$$

Contrastive KL Regularization

$$\begin{aligned} R(\pi_\psi; \pi^+, \pi^-) &:= -\bar{D}_{KL}(\pi^+ || \pi_\psi; \mathbf{s}_t^+, \mathbf{a}_t^+) + \bar{D}_{KL}(\pi^- || \pi_\psi; \mathbf{s}_t^-, \mathbf{a}_t^-) \\ &\approx \frac{1}{L} \sum_{l=1}^L \left[\log \frac{\pi^+(\mathbf{a}_{t+l}^+ | \mathbf{s}_{t+l}^+)}{\pi_\psi(\mathbf{a}_{t+l}^+ | \mathbf{s}_{t+l}^+)} - \log \frac{\pi^-(\mathbf{a}_{t+l}^- | \mathbf{s}_{t+l}^-)}{\pi_\psi(\mathbf{a}_{t+l}^- | \mathbf{s}_{t+l}^-)} \right] \end{aligned}$$

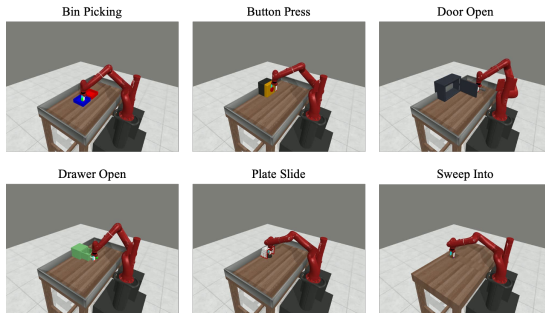
- Encourages π_ψ to align with preferred policy π^+ and diverge from less preferred π^- .
- Approximates into L -horizon undiscounted sum with sampled segments $\{\mathbf{s}_t^+, \mathbf{a}_t^+\} \sim \zeta^+$ and $\{\mathbf{s}_t^-, \mathbf{a}_t^-\} \sim \zeta^-$
- Mitigates likelihood mismatch over sequential rollouts.

Deterministic Pseudo Labeling

$$S_{\text{PPL-d}}(\pi_{\psi}; \zeta^+) - S_{\text{PPL-d}}(\pi_{\psi}; \zeta^-) = \sum_{t \geq 0} \left[\log \frac{\pi_{\psi}(\mathbf{a}_t^+ | \mathbf{s}_t^+)}{\pi_{\psi}(\mathbf{a}_t^- | \mathbf{s}_t^-)} + \frac{1}{L} \sum_{l=1}^L \log \frac{\pi_{\psi}(\mathbf{a}_{t+l}^+ | \mathbf{s}_{t+l}^+)}{\pi_{\psi}(\mathbf{a}_{t+l}^- | \mathbf{s}_{t+l}^-)} \right]$$

- Behavior policy is typically unknown in offline setting.
- Assign a pseudo label as if each segment is generated by deterministic policy.

Experiments



- 1. Is PPL robust when learning from heterogeneous datasets that include suboptimal data?
- 2. Does incorporating policy labels lead to improved performance?
- 3. Can PPL be applied effectively in an online setting?

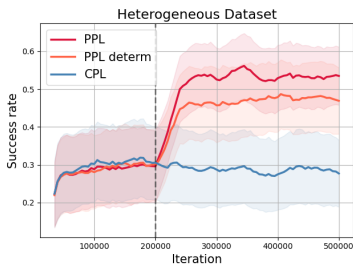
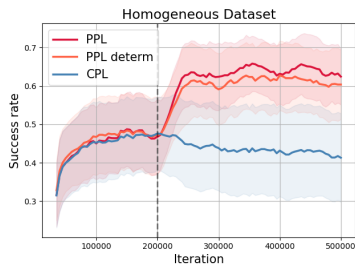
Offline MetaWorld Results

Table 2: Success rates of all methods across six tasks on the MetaWorld benchmark on different datasets. Each score is reported with the maximum average performance across four seeds over 200 episode evaluation window.

		Bin Picking	Button Press	Door Open	Drawer Open	Plate Slide	Sweep Into
Homogeneous Dense	SFT	39.7 \pm 19.2	71.5 \pm 3.3	48.0 \pm 15.6	56.2 \pm 1.8	64.8 \pm 0.8	70.0 \pm 6.5
	P-IQL	62.0 \pm 4.4	72.3 \pm 1.0	47.7 \pm 5.1	58.0 \pm 5.7	70.5 \pm 6.1	65.8 \pm 1.3
	CPL	22.7 \pm 5.5	64.3 \pm 1.4	29.0 \pm 4.3	54.0 \pm 4.3	65.5 \pm 3.1	69.8 \pm 3.3
	PPL	83.5 \pm 4.4	79.8 \pm 4.8	39.3 \pm 2.0	69.2 \pm 5.5	64.7 \pm 2.0	72.8 \pm 4.8
Homogenous Sparse	SFT	33.5 \pm 5.4	67.4 \pm 1.5	31.3 \pm 2.1	54.9 \pm 2.7	67.1 \pm 3.7	78.3 \pm 2.5
	P-IQL	72.4 \pm 6.6	74.5 \pm 0.0	58.5 \pm 1.4	51.4 \pm 4.6	76.3 \pm 1.6	79.0 \pm 2.6
	CPL	26.5 \pm 1.0	63.7 \pm 1.3	28.5 \pm 5.8	50.1 \pm 4.5	65.1 \pm 2.8	72.9 \pm 6.1
	PPL	87.2 \pm 3.5	87.3 \pm 2.8	49.3 \pm 6.5	68.5 \pm 5.3	64.0 \pm 6.4	73.9 \pm 3.5
Heterogeneous Dense	SFT	18.5 \pm 23.8	63.7 \pm 12.2	26.0 \pm 12.5	32.0 \pm 5.7	62.8 \pm 1.6	53.0 \pm 9.1
	P-IQL	51.2 \pm 5.3	62.5 \pm 4.9	32.0 \pm 3.5	41.8 \pm 3.8	67.0 \pm 3.0	59.3 \pm 3.7
	CPL	1.2 \pm 0.8	49.7 \pm 3.0	17.3 \pm 2.5	26.0 \pm 2.2	59.2 \pm 7.7	51.2 \pm 3.0
	PPL	59.7 \pm 18.6	73.8 \pm 3.3	25.8 \pm 2.0	58.5 \pm 3.8	69.8 \pm 2.3	57.3 \pm 8.6
Heterogeneous Sparse	SFT	12.2 \pm 1.0	63.7 \pm 4.7	17.8 \pm 0.8	38.7 \pm 3.0	70.7 \pm 3.8	60.7 \pm 2.5
	P-IQL	48.0 \pm 5.6	71.0 \pm 6.6	44.1 \pm 3.2	47.5 \pm 3.0	72.0 \pm 4.0	64.3 \pm 1.0
	CPL	18.0 \pm 6.1	50.8 \pm 0.8	18.5 \pm 3.0	32.1 \pm 1.6	67.3 \pm 5.5	55.5 \pm 3.3
	PPL	83.8 \pm 3.8	83.5 \pm 1.8	34.3 \pm 7.6	60.8 \pm 7.3	71.2 \pm 1.9	63.3 \pm 4.2

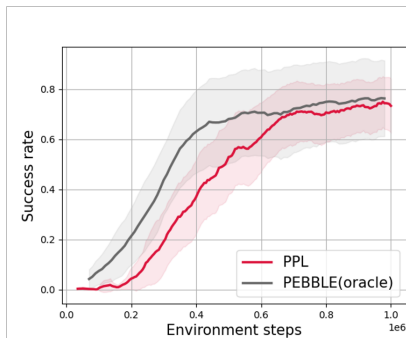
- PPL outperforms CPL and P-IQL especially in sparse and heterogeneous settings.
- Robust across 6 MetaWorld tasks.

Ablation on Policy Labels



- Deterministic pseudo-labels perform worse in heterogeneous data.
- Shows benefit of incorporating true or approximated behavior policies.

Online RLHF Setting



Method	Params
PPL	76k
PEBBLE	859k

- PPL can be directly applied to online setting.
- Achieves competitive performance with 1/10 of PEBBLE's parameters.

Conclusion

- PPL resolves likelihood mismatch by modeling regret w.r.t. behavior policies.
- Theoretical foundations show regret minimization \Leftrightarrow forward KL.
- Contrastive KL regularization provides robustness across offline and online RLHF.
- PPL is sample-efficient and scalable for real-world RLHF tasks.

Thank you!

Project Page



References



[Gaon An \(2023\)](#)

Direct preference-based policy optimization without reward modeling



[Paul F Christiano \(2017\)](#)

Deep reinforcement learning from human preferences



[Joey Hejna \(2024\)](#)

Contrastive preference learning: Learning from human feedback without reinforcement learning



[W. Bradley Knox \(2024\)](#)

Learning optimal advantage from preferences and mistaking it for reward.



[Kimin Lee \(2021\)](#)

Pebble: Feedback-efficient interactive reinforcement learning via relabeling experience and unsupervised pre-training.

References



Rafael Rafailov (2024)

From r to q^* : Your language model is secretly a q-function.



Rafael Rafailov (2024)

Direct preference optimization: Your language model is secretly a reward model.



Omar Shaikh (2024)

Show, Don't Tell: Aligning Language Models with Demonstrated Feedback.



Yongcheng Zeng (2024)

Token-level Direct Preference Optimization



Brian D. Ziebart (2010)

Modeling purposeful adaptive behavior with the principle of maximum causal entropy.