

FUSING REWARDS AND PREFERENCES IN REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

We present Dual-Feedback Actor (DFA), a reinforcement learning algorithm that fuses both individual rewards and pairwise preferences (if available) into a single update rule. DFA uses the policy’s log-probabilities directly to model the preference probability, avoiding a separate reward-modeling step. Preferences can be provided by human-annotators (at state-level or trajectory-level) or be synthesized online from Q-values stored in an off-policy replay buffer. Under a Bradley–Terry model, we prove that minimizing DFA’s preference loss recovers the entropy-regularized Soft Actor-Critic (SAC) policy. Our simulation results show that DFA trained on generated preferences matches or exceeds SAC on six control environments and demonstrates a more stable training process. With only a semi-synthetic preference dataset under Bradley-Terry model, our algorithm outperforms reward-modeling reinforcement learning from human feedback (RLHF) baselines in a stochastic GridWorld and approaches the performance of an oracle with true rewards.

1 INTRODUCTION

Over the past decade, Reinforcement Learning (RL) has achieved remarkable success across a wide range of applications, including video games (Knox & Stone, 2008; Warnell et al., 2018), recommendation systems (Kohli et al., 2013; Zeng et al., 2016), and autonomous driving (Kiran et al., 2021). RL focuses on how agents make decisions while interacting with dynamic, changing environments. At each time step, an agent chooses an action based on its current state and receives a reward that indicates how good that action was. The goal is to learn a policy that maximizes the total reward accumulated over time. In traditional RL, the reward function is usually manually designed by experts to guide the agent’s behavior toward desired outcomes. However, crafting such a function is a challenging and often ambiguous task (Ng et al., 2000).

To overcome the limitations of hand-engineered rewards, Reinforcement Learning from Human Feedback (RLHF) has emerged as a compelling alternative, particularly in the fine-tuning of large language models (LLMs) (Christiano et al., 2017; Stiennon et al., 2020; Ouyang et al., 2022). RLHF bypasses manual reward specification by inferring a reward model from human preferences over trajectory pairs. This reward model then guides policy optimization using standard RL algorithms. Despite its empirical successes, RLHF methods relying on reward inference, face significant practical and theoretical challenges, including reward model misspecification, overfitting, distribution shift, and non-identifiability of reward functions (Zhu et al., 2024; Casper et al., 2023). Moreover, the reward inference step introduces additional complexity and often requires large volumes of annotated data.

To simplify the pipeline and avoid reward inference, in the context of language modeling, Direct Preference Optimization (DPO) has recently been proposed as a direct approach to exploit human preferences (Rafailov et al., 2023). Thanks to a closed-form expression of the optimal policy under a Bradley–Terry preference model, DPO avoids estimating the reward function. Although DPO has shown promising results in fine-tuning large language models, its loss formulation tends to induce deterministic policies and is susceptible to mode collapse (Azar et al., 2024; Sharifnassab et al., 2024). Moreover, the existing theory for DPO only covers contextual bandits or MDPs with deterministic transitions (Rafailov et al., 2023; 2024). As a result, directly applying DPO (or methods suggested in Guo et al. (2024); Xie et al. (2024)) in general reinforcement learning settings is suboptimal, where effective exploration is critical for policy improvement in stochastic MDPs (Zhang & Ying, 2024).

054 More recently, ZPG (Zhang & Ying, 2024) suggested an RLHF approach that does not rely on a
055 reward model and is designed for non-deterministic MDPs. However, as the authors acknowledged,
056 the algorithm lacks a strategic exploration mechanism. Furthermore, it relies on trajectory-level
057 preference comparisons and performs on-policy updates, hence previously collected data are not
058 reused.

059 In this work, we introduce Dual-Feedback Actor (DFA), a reinforcement learning algorithm that works
060 for stochastic MDPs and unifies scalar rewards and preference-based feedback into a single, principled
061 policy update rule. Unlike many prior approaches in RLHF that infer a separate reward model from
062 preferences, DFA directly incorporates preferences into the policy optimization objective using the
063 policy’s log-probabilities and retains Soft Actor-Critic (SAC)-style entropy-driven exploration. The
064 main contributions are as follows:

- 065 • Our approach offers dual compatibility with both rewards and preferences. When numerical
066 rewards are available, the agent updates its Q-networks and incorporates preference-based
067 learning by synthesizing preferences from Q-values. This dual approach allows the agent to
068 use reward signals while maintaining flexibility to incorporate human feedback, especially
069 in settings where rewards are sparse or absent.
- 070 • Our approach can be used not only in on-policy manner but also in off-policy manner, which
071 enables more sample-efficient learning by reusing past experiences stored in a replay buffer.
072 This is particularly valuable for hierarchical RL applications where sample efficiency is
073 needed to train the policy of each layer.
- 074 • Under the assumptions stated in Section 5, we prove that minimizing DFA’s preference loss
075 recovers the entropy-regularized SAC solution, formally bridging preference optimization
076 and entropy-regularized RL. Consequently, DFA inherits SAC’s entropy-driven exploration,
077 maintaining diverse action sampling even when it learns solely from preferences.
- 078 • Experimental results in Section 6 show that DFA consistently matches or outperforms both
079 reward-based and preference-based baselines on six control tasks and a stochastic GridWorld,
080 while yielding a more stable training process.

083 2 RELATED WORK

084
085 There are two dominant paradigms for incorporating human feedback in reinforcement learning. The
086 first relies on reward modeling: These methods first fit a scalar reward (or value) prediction model
087 from preference data and then treat this learned reward model as the surrogate reward for standard
088 policy optimization. This two-stage pipeline was introduced in (Christiano et al., 2017), Schoenauer
089 et al. (2014), and later scaled to large language models by Ziegler et al. (2019), Stiennon et al. (2020),
090 and Ouyang et al. (2022). The second relies on direct policy optimization: These algorithms bypass
091 an explicit reward model and update the policy parameters solely from preference comparisons
092 (Wilson et al., 2012; Busa-Fekete et al., 2014; Akrouf et al., 2011).

093 For reward-modeling approaches, several works (Saha et al., 2023; Zhu et al., 2023; Wu & Sun, 2023)
094 consider linearly parameterized reward models and characterize the error bounds of the estimated
095 parameters, and prove that subsequent reward-based RL can tolerate small errors in rewards. Zhan
096 et al. (2023) extend this analysis to more general reward function classes under some conditions.
097 These analyses have been extended to direct policy optimization approaches in Xu et al. (2020); Chen
098 et al. (2022); Zhang et al. (2024).

099 In the context of language modeling, DPO (Rafailov et al., 2023) provides a direct approach to
100 aligning language models with human preferences by optimizing a policy to maximize the likelihood
101 of preferred responses over nonpreferred ones, eliminating the need for an explicit reward model.
102 SPO (Sharifnassab et al., 2024) optimizes model output directly over a preference dataset through the
103 natural conditional probability of the preferred responses over nonpreferred ones. Similar approaches
104 have also been explored in this literature (Xu et al., 2024; Ethayarajh et al., 2024; Hong et al., 2024;
105 Park et al., 2024; Hong et al., 2024; Meng et al., 2024; Li et al., 2025). RLHF has also been studied
106 in other aspects. For example, the framework in Swamy et al. (2024) casts RLHF as a two-player
107 zero-sum game. However, they still estimate the rewards (and subsequently apply PPO, TRPO, or
SAC) based on a constantly updated queue of recent rollouts, which can cause data staleness issues.

Recent work, Xie et al. (2024), inspired by DPO, combines DPO with optimistic exploration to design XPO in the function approximation regime with provable convergence. ZPG (Zhang & Ying, 2024) aims to address RLHF without relying on a reward model and is designed for non-deterministic MDPs. However, it lacks an exploration mechanism, which is essential for general RL applications. Although previous work brought advancement in several aspects, existing algorithms are rarely benchmarked (theoretically and experimentally) against strong reward-based baselines such as SAC.

3 PRELIMINARIES

In this section, we introduce the notation for RL, RLHF, and review the DPO objective (Rafailov et al., 2023). We model the environment as a finite-horizon Markov Decision Process (MDP). An MDP can be represented as a tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma, p_0 \rangle$, where \mathcal{S} and \mathcal{A} are state space and action space, respectively. The conditional probability of transition from state s to s' with action a is denoted by $P(s'|s, a)$. The probability distribution over the initial state s_0 is denoted by $p_0(s_0)$. The parameter $\gamma \in (0, 1)$ denotes the discount factor. At each time step t , $r(s_t, a_t)$ returns the reward of taking action a_t in the state s_t . Actions are chosen according to the policy π where $\pi(a|s)$ is the probability of taking action a for a given state s . Here, we assume that the policy is parameterized with a vector $\theta \in \mathbb{R}^d$ and use shorthand notation π_θ for $\pi_\theta(a|s)$. For a given time horizon H , we define $\tau = (s_0, a_0, \dots, s_{H-1}, a_{H-1})$ as a sequence of state-action pairs called a trajectory. $R(\tau)$ is a function that returns the discounted accumulated reward of each trajectory as follows: $R(\tau) := \sum_{h=0}^{H-1} \gamma^h r(s_h, a_h)$ where $\gamma \in (0, 1)$ is the discount factor.

Given a policy π , the *state-value function* and the *action-value function* (or *Q-function*) are

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \mid s_0 = s \right], Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

Classical RLHF feedback setting (Christiano et al., 2017). Let $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma, p_0 \rangle$ be the finite-horizon MDP where the true reward $r(s, a)$ is *hidden*. Hence, we ask humans to compare trajectories and form the preference dataset

$$\mathcal{D} = \{(\tau_k^+, \tau_k^-)\}_{k=1}^K, \quad \tau_k^+ \succ \tau_k^-,$$

where K is the total number of pairs, and τ_k^+ is *preferred* to τ_k^- which is denoted as $\tau_k^+ \succ \tau_k^-$. We define a parametric function $r_\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ to *approximate* the latent reward. For any trajectory $\tau = (s_0, a_0, \dots, s_{H-1}, a_{H-1})$, we define the model return as:

$$R_\phi(\tau) = \sum_{h=0}^{H-1} \gamma^h r_\phi(s_h, a_h).$$

The parameters ϕ are learned by maximum likelihood under the Bradley–Terry model (Bradley & Terry, 1952), which is equivalent to minimizing the following loss:

$$\mathcal{L}(\phi) = - \mathbb{E}_{(\tau^+, \tau^-) \sim \mathcal{D}} [\log \sigma(R_\phi(\tau^+) - R_\phi(\tau^-))], \quad \sigma(z) = \frac{1}{1 + e^{-z}}.$$

Let \hat{r}_ϕ be the estimated reward function. Next, with \hat{r}_ϕ fixed, a policy-gradient method such as PPO (Schulman et al., 2017) or SAC (Haarnoja et al., 2018) updates π_θ to maximize

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [\hat{R}_\phi(\tau)].$$

A well-known drawback on this two-stage pipeline is its sensitivity to noise, and any overfitting in \hat{r}_ϕ propagates directly to the final policy updates (Casper et al., 2023).

DPO in language models (Rafailov et al., 2023) In language models, a *state* is the text prefix (or prompt) x , and an *action* is the response y produced by the model (call it continuations). Annotators make a choice among two full continuations (y^+, y^-) sampled from the same prompt, giving the preference dataset

$$\mathcal{D} = \{(x_k, y_k^+, y_k^-)\}_{k=1}^K, \quad y_k^+ \succ y_k^-.$$

Let π_{ref} be the frozen base model (e.g. a pre-trained GPT checkpoint). For a prompt x and two candidate continuations y^+, y^- , define the log-probability gap:

$$\Delta_{x,y^+,y^-}(\theta) = [\log \pi_{\theta}(y^+ | x) - \log \pi_{\theta}(y^- | x)] - [\log \pi_{\text{ref}}(y^+ | x) - \log \pi_{\text{ref}}(y^- | x)].$$

$\Delta_{x,y^+,y^-}(\theta)$ captures how much more the new model prefers the chosen continuation over the rejected one, *relative* to the base model. DPO then minimizes

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(x,y^+,y^-) \sim \mathcal{D}} \left[\log \sigma(\alpha \Delta_{x,y^+,y^-}(\theta)) \right], \quad \sigma(z) = \frac{1}{1+e^{-z}}, \quad \alpha > 0.$$

Minimizing \mathcal{L}_{DPO} pushes the new model toward the preferred continuation, while limiting it to the safe behavior of π_{ref} . The absence of a separate reward model in DPO removes a major source of overfitting or noisy evaluations of the reward modeling. However, DPO assumes a Bradley-Terry choice model to derive its loss function, and this loss tends to produce near-deterministic models. This reduced diversity makes DPO prone to mode collapse (Azar et al., 2024).

4 METHODS

In this section, we introduce our *Dual-Feedback Actor* (DFA). In order to describe the DFA algorithm, we first introduce the state-wise feedback setting as follows:

State-wise feedback. In this setting, the annotator does *not* compare full trajectories. Instead, at a given state s_k , the annotator sees two actions, marks the winner a_k^+ over the loser a_k^- . Then, the following preference dataset is formed:

$$\mathcal{D} = \{(s_k, a_k^+, a_k^-)\}_{k=1}^K, \quad a_k^+ \succ a_k^-,$$

where a^+ is *preferred* to a^- at state s_k . In the subsections below, we first consider the case where the agent learns only from state-wise human comparisons. Second, we show how to synthesize preferences from numerical rewards when they are available. Finally, we extend DFA to trajectory-based comparisons.

4.1 LEARNING WITH ONLY STATE-WISE PREFERENCES

Assume we have collected a set of preference comparisons

$$\mathcal{D} = \{(s_k, a_k^+, a_k^-)\}_{k=1}^K, \quad a_k^+ \succ a_k^-.$$

Unlike classical RLHF, we do not assume an underlying Bradley-Terry reward model. Instead, we rely on the policy’s log-probabilities to model the preference probability directly. For any pair (s, a^+, a^-) we define the *preference probability* produced by the current policy π_{θ} as

$$P_{\theta}(a^+ \succ a^- | s) = \frac{\pi_{\theta}(a^+ | s)^{\alpha}}{\pi_{\theta}(a^+ | s)^{\alpha} + \pi_{\theta}(a^- | s)^{\alpha}}, \quad \alpha > 0. \quad (1)$$

The exponent α controls the uncertainty assigned to the policy’s output: $\alpha \rightarrow 0$ yields a nearly uniform (high-entropy) choice, while $\alpha \rightarrow \infty$ approaches a hard winner-takes-all rule.

The negative log-likelihood equation 1 gives the state-wise preference loss:

$$\mathcal{L}_{\text{pref}}(\theta) = -\mathbb{E}_{(s,a^+,a^-) \sim \mathcal{D}} [\log P_{\theta}(a^+ \succ a^- | s)]. \quad (2)$$

Minimizing $\mathcal{L}_{\text{pref}}$ directly increases the probability that π_{θ} selects the human-preferred action, without introducing auxiliary reward networks or relying on any latent utility assumptions¹. Note that equation 2 can be reformulated as follows:

¹Eq. equation 2 is identical to the *preference loss* $\mathcal{L}_{\text{pref}}^{\alpha}$ used in Soft Preference Optimization (SPO) (Sharifnassab et al., 2024). Although DFA adopts the same logistic pairwise-loss form, the similarity ends there. In SPO, the same term is combined with a global KL regularizer $D_{\text{KL}}(\pi_{\theta} \| \pi_{\text{ref}})$, whereas here we study the stand-alone preference part and show that, under some assumptions, it aligns the policy with the entropy-regularized RL solution (Theorem 5.2). Moreover, SPO is in the context of LLMs and is designed for an offline setting. DFA targets stochastic MDPs, supports off-policy replay, preserves SAC-style entropy exploration with theoretical analysis, and unifies numeric rewards with preferences. Synthesizing preferences, as explained in Section 4.2, is another key innovation in DFA that allows for online settings in RL.

$$\mathcal{L}_{\text{pref}}(\theta) = -\mathbb{E}_{(s, a^+, a^-) \sim \mathcal{D}} \left[\log \sigma \left(\alpha \left(\log \pi_{\theta}(a^+ | s) - \log \pi_{\theta}(a^- | s) \right) \right) \right]$$

In simulated environments or settings where numerical rewards are accessible, it is possible to synthesize preference data from these rewards or their proxies, such as Q-values. Our approach, introduced in the next section, is particularly useful when integrating preference-based learning into an agent’s training loop, even when direct human feedback is unavailable or insufficient. Our method fuses numerical rewards and preference data by synthesizing preferences from numerical rewards.

4.2 SYNTHESIZING PREFERENCES FROM NUMERICAL REWARDS

We use Q-values as a proxy to create preference pairs, enabling online preference generation during policy updates without explicitly constructing full trajectory segments. Estimating Q-values can be done through any method in the literature, and is particularly relevant in off-policy methods such as SAC, where a replay buffer stores past experiences as tuples (s_t, a_t, r_t) , where s_t , a_t , and r_t are state, action and reward at time t , respectively.

Our approach works as follows: For a batch of states $\{s_i\}_{i=1}^N$ sampled from the replay buffer, we generate two candidate actions to form preference pairs: The first action, denoted by a_i , corresponds to the action originally taken in state s_i as stored in the replay buffer. This action reflects the historical behavior of the agent at the time the state was visited. The second action, denoted by a'_i , is obtained from the replay buffer by identifying the action associated with the nearest state to s_i (denote it with s'_i)². For both actions, we compute their respective Q-values. The action with the higher Q-value is designated as the preferred action a_i^+ , while the other is labeled as the rejected action a_i^- :

$$\begin{aligned} \text{If } Q(s_i, a_i) > Q(s_i, a'_i), \quad & \text{then } (a_i^+, a_i^-) = (a_i, a'_i), \\ & \text{else } (a_i^+, a_i^-) = (a'_i, a_i). \end{aligned}$$

This process effectively synthesizes preference data in the form of state-action pairs $\mathcal{D}^{\text{Syn}} = \{(s_i, a_i^+, a_i^-)\}_{i=1}^N$ for each batch. The loss in equation 2, is then calculated over the states $\{s_i\}_{i=1}^N$ and using their associated preferred and rejected actions. Specifically, the loss encourages the policy to assign higher probability to preferred actions over rejected ones, scaled by the parameter α

$$\mathcal{L}_{\text{pref}}^{\text{Syn}}(\theta) = -\mathbb{E}_{(s_i, a_i^+, a_i^-) \sim \mathcal{D}^{\text{Syn}}} \left[\log \left(\sigma \left(\alpha \left(\log \pi_{\theta}(a_i^+ | s_i) - \log \pi_{\theta}(a_i^- | s_i) \right) \right) \right) \right]$$

where $\sigma(\cdot)$ is the sigmoid function. Figure 1 gives a high-level schematic of our methodology.

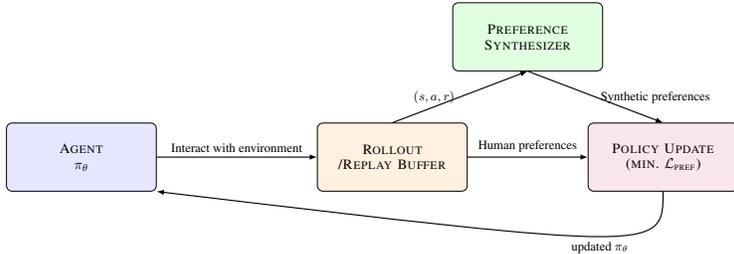


Figure 1: Data flow in DFA. The agent executes its current policy π_{θ} and stores the transitions (may include reward-based transitions or human-annotated preferences). If reward-based transitions are available, the *Preference Synthesizer* can convert them into synthetic preference pairs. This process can be done in either an on-policy or off-policy fashion. Both human and synthetic preferences can be used in *Policy-Update*, which minimizes the preference loss $\mathcal{L}_{\text{pref}}$ and outputs an improved policy.

²In our experiments, we compute the Euclidean distance between s_i and all states in the buffer, select the closest state s'_i , and retrieve its corresponding action a'_i .

4.3 EXTENSION OF THE LOSS TO TRAJECTORY-BASED COMPARISONS

State-wise comparisons can be easy to collect (for instance, a single frame rather than a full video in video games) and give richer training signals, but one may prefer to rank the whole trajectories (Christiano et al., 2017; Zhang & Ying, 2024), hence, we extend DFA to accept trajectory-level preferences as well. For trajectory-level comparisons, we store pairs

$$D^{\text{traj}} = \{(\tau_k^+, \tau_k^-)\}_{k=1}^K, \quad \tau = (s_1, a_1, \dots, s_T).$$

The policy assigns a likelihood to any full trajectory as follows: $\pi_\theta(\tau) = \prod_{t=1}^T \pi_\theta(a_t | s_t)$. The preference probability is the same as before, but now in terms of trajectory likelihoods:

$$P_\theta^{\text{traj}}(\tau^+ \succ \tau^-) = \frac{\pi_\theta(\tau^+)^\alpha}{\pi_\theta(\tau^+)^\alpha + \pi_\theta(\tau^-)^\alpha}, \quad \alpha > 0. \quad (3)$$

The negative log-likelihood of equation 3 gives trajectory-based preference loss:

$$\mathcal{L}_{\text{pref}}^{\text{traj}}(\theta) = -\mathbb{E}_{(\tau^+, \tau^-) \sim D^{\text{traj}}} \left[\log P_\theta^{\text{traj}}(\tau^+ \succ \tau^-) \right]. \quad (4)$$

5 THEORETICAL ANALYSIS

In this section, we show that, under Bradley–Terry model on the soft optimal Q -function, minimizing our preference loss is equivalent to recovering the optimal policy for entropy-regularized reinforcement learning (Haarnoja et al., 2017). Concretely, we analyze the tabular setting for the state-wise preferences and identify its unique minimizer. This establishes the equivalence of preference optimization and entropy-regularized RL. We should emphasize that the BT model is not a requirement of DFA Algorithm; it is only used to derive Theorem 5.2 in the following. A trajectory-wise analysis and its connection to the state-wise analysis, is provided in Appendix B.

Assumption 5.1 (Bradley–Terry preferences on the *soft*-optimal Q -function). Let $Q^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be the soft-optimal state-action value function of the MDP, i.e.,

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \mathcal{H}(\pi(\cdot | s_t))) \mid s_0 = s, a_0 = a \right],$$

where $\mathcal{H}(\cdot)$ is the entropy function and λ is the entropy coefficient. Assume that there exists a parameter $\beta > 0$ such that, for every $s \in \mathcal{S}$ and any $a, b \in \mathcal{A}$,

$$P^*(a \succ b | s) = \sigma\left(\beta [Q^*(s, a) - Q^*(s, b)]\right), \quad \sigma(z) = \frac{1}{1 + e^{-z}}.$$

Theorem 5.2 (Preference loss recovers the optimal policy). *Fix a state $s \in \mathcal{S}$ and abbreviate $Q_a^* := Q^*(s, a)$. Suppose that Theorem 5.1 holds. Under uniform sampling of ordered pairs $(a, b) \sim \text{Unif}(\mathcal{A}^2)$ and the tabular full-support parameterization $\ell_a = \log \pi(a | s)$ ($\sum_a e^{\ell_a} = 1$, $e^{\ell_a} > 0$), consider the preference loss*

$$\mathcal{L}(\ell) = -\frac{1}{|\mathcal{A}|^2} \sum_{(a,b) \in \mathcal{A}} P^*(a \succ b | s) \log \sigma(\alpha(\ell_a - \ell_b)), \quad \alpha > 0.$$

This loss is strictly convex on the set of full-support policies and is minimized uniquely at

$$\pi_\star(a | s) = \frac{\exp\left(\frac{\beta}{\alpha} Q_a^*\right)}{\sum_{a' \in \mathcal{A}} \exp\left(\frac{\beta}{\alpha} Q_{a'}^*\right)}. \quad (5)$$

Furthermore, this Gibbs distribution coincides with the global maximizer of the entropy-regularized RL (or SAC objective) when $\lambda = \alpha/\beta$.³

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \mathcal{H}(\pi(\cdot | s_t))) \right]. \quad (6)$$

³All proofs are provided in the Appendix.

Theorem 5.2 states that, when human (or synthetic) comparisons follow a Bradley-Terry model whose latent utility equals the ground truth Q^* , the preference loss is perfectly aligned with the entropy-regularized control objective (Haarnoja et al., 2017). The optimizer equation 5 is soft-max policy whose inverse temperature is the ratio β/α : The parameter β captures how consistently the annotator prefers higher-value actions, while the parameter α adjusts the learner’s uncertainty. In particular, setting $\lambda = \alpha/\beta$ recovers the SAC trade-off between exploitation (large β) and exploration (large α) (Haarnoja et al., 2018).

Remark 5.3. If the Bradley-Terry assumption holds for any *arbitrary* soft state-action value function, for instance, the current critic estimate $Q_k(s, a)$ in SAC (Haarnoja et al., 2018). Then Theorem 5.2 implies that the preference loss is minimized by

$$\pi_{k+1}(a | s) = \frac{\exp(\frac{\beta}{\alpha} Q_k(s, a))}{\sum_{a'} \exp(\frac{\beta}{\alpha} Q_k(s, a'))}.$$

This update is *exactly* the policy-improvement step in SAC that maximizes the entropy-regularized objective

$$\max_{\pi} \left\{ \mathbb{E}_{a \sim \pi} [Q_k(s, a)] + \lambda \mathcal{H}(\pi(\cdot | s)) \right\}.$$

Hence, as the critic converges ($Q_k \rightarrow Q^*$), repeated minimization of the preference loss yields the soft-optimal SAC policy. Therefore, preference learning can be viewed as performing policy improvement in SAC, but driven solely by comparative feedback.

Remark 5.4. The assumption that $(a, b) \sim \text{Unif}(\mathcal{A}^2)$ in Theorem 5.2 is made for simplicity of analysis. In practice, one can approximate this condition by drawing a mini-batch of ordered pairs at each update and down-sampling (or re-weighting) each pair by the inverse of its frequency in the batch; This produces a uniform sub-sampled action pairs required by the theorem.

6 EXPERIMENTAL RESULTS

In this section, we benchmark DFA against prior work. The complete code is provided in the supplementary material. We first compare DFA with the reward-based baseline SAC (hence, we have to synthesize preferences following Section 4.2), and then against recent preference-based methods.

6.1 COMPARISON WITH SAC VIA SYNTHETIC PREFERENCES

In this section, we evaluate the proposed algorithm (DFA) and compare it with related work on six control tasks in MuJoCo (Todorov et al., 2012), a physics simulator known for fast and accurate simulations in areas such as robotics, biomechanics, and graphics. Since published benchmarks (e.g., OpenAI SpinningUp) consistently identify SAC as the strongest baseline on many environments, we compare DFA exclusively with SAC. We briefly explain the six environments we consider in Appendix C.

In Figure 2, we monitor the average episode return versus system probes, which represents the total number of environment interactions. In this experiment, DFA continually generates synthetic preference pairs from numerical rewards following 4.2. The underlying RL settings and replay buffer are identical to those of SAC. Both DFA and SAC run for 10×10^6 system probes across 5 different random seeds. We use mini-batch size 256 for both algorithms. A new preference batch of size $N = 256$ is created during every gradient step. Figure 2 shows that DFA matches or exceeds SAC on Walker2d, Hopper, Swimmer, and Humanoid. In MountainCarContinuous, we could not find SAC settings that produced learning, a problem others have reported as well.⁴ DFA, in contrast, learned a good policy on this task with the same range of hyperparameters used for the other environments.

Interestingly, DFA’s learning curves are noticeably smoother, while SAC exhibits significant fluctuations. We attribute this stability to the synthesized preference pairs, which are constructed according to Section 4.2, and appear to act as an implicit denoising regularizer. We note that reducing the learning rate or adjusting other hyperparameters to avoid fluctuations for SAC resulted in lower average returns, thus, we maintained the higher learning rate configuration to ensure fair comparison.

⁴<https://github.com/rail-berkeley/softlearning/issues/76>

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

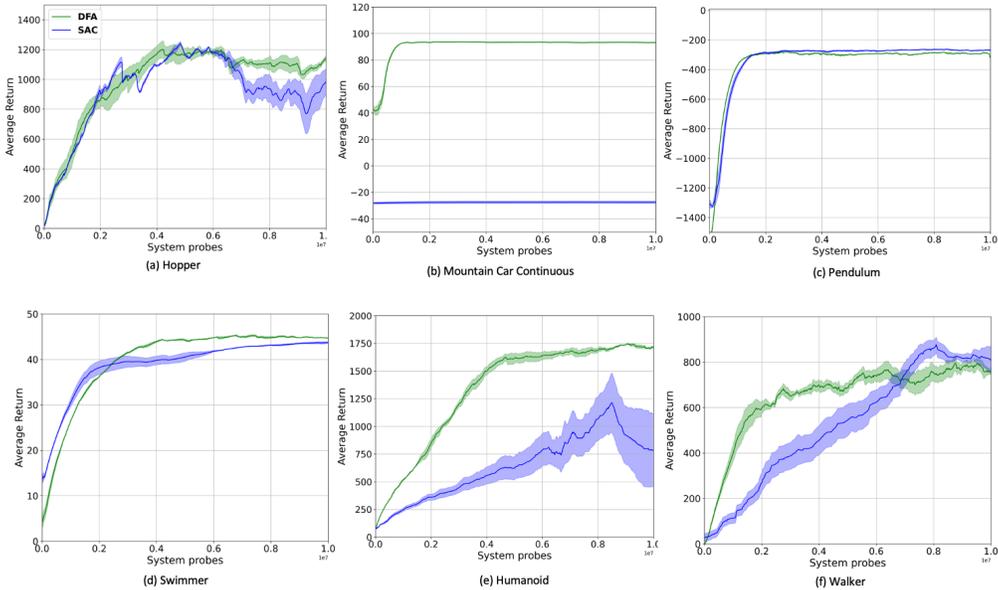


Figure 2: DFA (green) vs. SAC (blue) on the six MuJoCo control tasks. DFA matches or exceeds SAC and shows smoother training. The solid line is the mean episode return, and the shaded region shows an 90% confidence interval over 5 seeds.

These results confirm the claim of Theorem 5.2: *once we use preference data aligned with the optimal Q-values, numerical rewards can be dropped without losing performance*. This unifies reward-free human alignment and reward-based RL under a single log-likelihood objective.

6.2 COMPARISON WITH RM METHODS

In this section, we evaluate our DFA algorithm against traditional reward modeling approaches in the context of learning from human preferences. While the previous section demonstrated DFA’s effectiveness with generated preferences derived from numerical rewards, here we focus on the more challenging scenario where only human comparative feedback is available, without access to ground-truth rewards.

We conduct experiments in a stochastic GridWorld environment, which provides a controlled testbed for preference-based learning (Zhang & Ying, 2024). In this environment, the agent starts at the center of the grid and can take four actions: up, down, left, or right. The environment includes the following aspects: (1) To build the ground play, a coin is flipped for each cell, and if heads, a reward sampled from $\mathcal{N}(0, 1)$ is placed in that cell; (2) While the agent is moving, with probability 0.4, the chosen action is reversed (e.g., "up" becomes "down"). Each episode has a fixed horizon of 20 steps, and the agent’s goal is to maximize the cumulative reward collected. This environment is particularly suitable for preference-based learning evaluation as it combines stochastic dynamics with a non-trivial reward structure that requires exploration.

To simulate human preferences, we simulate a panel of annotators who provide comparative feedback between trajectories. Following standard practice in RLHF literature, we model the annotator’s preference probability using the Bradley-Terry model: $P(\tau_1 \succ \tau_0) = \sigma(R_1 - R_0)$, where R_i is the cumulative reward of trajectory τ_i and σ is the sigmoid function. For robustness, each preference query aggregates votes from M independent annotators, and the majority vote determines the final preference. This approach simulates the noise and variability in real human feedback while maintaining a consistent underlying reward structure. For more implementation details, please see Appendix C. We compare DFA against the following approaches: **RM+PPO**: A two-stage approach that first learns a reward model from preference data using maximum likelihood estimation, then optimizes a policy using Proximal Policy Optimization (PPO) with the learned reward function. **ZPG (Zhang & Ying, 2024)**: A state-of-the-art RLHF method which estimates the policy gradient from

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

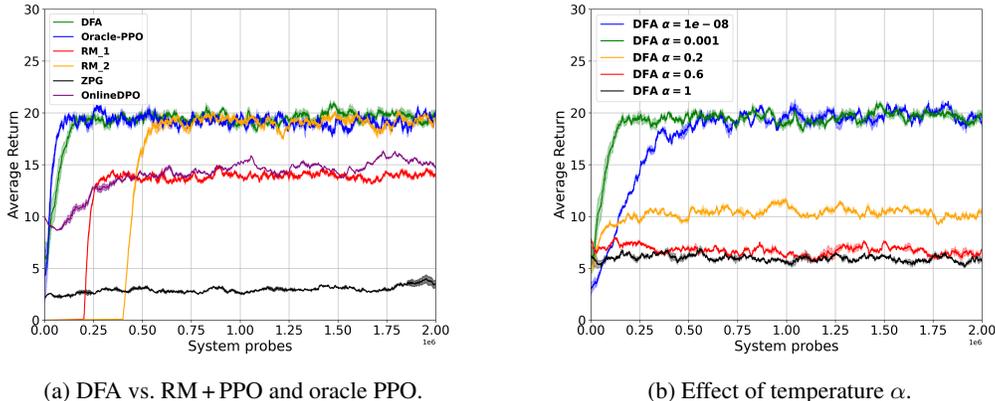


Figure 3: GridWorld results. (a) DFA learns faster and achieves higher rewards than reward-modeling baselines, approaching the oracle that has access to the true reward. (b) Effect of the temperature parameter α : a small but not too small value balances exploration and exploitation. Shaded regions denote 90% confidence intervals across 5 random seeds.

preference differences without fitting a reward model. **Oracle-PPO (upper bound)**: PPO directly on the *true* MDP reward r (it is unavailable in practice, but gives an upper bound on the performance.). **OnlineDPO**: We also include the recently-proposed OnlineDPO algorithm (Guo et al., 2024) as a direct-preference baseline.

Figure 3a demonstrates that DFA consistently outperforms reward modeling methods and performs comparably to Oracle-PPO, which has access to the true reward function. In this experiment, we compare against two variants of RM+PPO: RM_1 uses 200k environment steps for training the reward model, while RM_2 uses twice as many samples (400k steps). Despite the increased data budget for RM_2, DFA is still converging faster, highlighting the benefits of avoiding the two-stage pipeline. For ZPG, we used the official repository. Despite our efforts (and implementation tricks such as normalized gradient and gradient clipping), it could not be tuned to outperform the results shown in Figure 3a. We also re-implemented ZPG directly from the paper, closely matching the reported hyperparameters and implementation choices; we report the best observed performance across both implementations. In Figure 3a we use an annotator pool of $M = 500$; runs with smaller M and more complex environments show the same pattern and are included in Appendix C.

Figure 3b highlights the sensitivity of DFA to the parameter α . As shown in Figure 3b, setting α too high ($\alpha = 1.0$) gives almost no learning signal, while moderate values in the range 0.2–0.6 yield better results. The best result comes at $\alpha = 0.001$. When α is pushed to very small values (e.g., 10^{-8}), performance drops again because the policy becomes overly stochastic. These results suggest that α should be small but not too small to balance exploration and exploitation.

7 CONCLUSION

Dual-Feedback Actor (DFA) unifies scalar rewards and pairwise preferences in a single loss; when preferences follow a Bradley–Terry model on the optimal soft Q -function, this loss recovers the entropy-regularized SAC solution, formally linking reward- and preference-based RL. Empirically, DFA matches or exceeds SAC and outperforms reward-modeling baselines while training more smoothly. The main limitations are the Bradley–Terry assumption, the noise inherited by synthetic preferences from early Q estimates, and the computational cost of finding the nearest state in the replay buffer. The future work can be investigating other assumptions and evaluating DFA on larger, real human-in-the-loop tasks.

REFERENCES

- 486
487
488 P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Optimization algorithms on matrix manifolds.
489 In *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- 490
491 Riad Akrou, Marc Schoenauer, and Michele Sebag. Preference-based policy learning. In *Machine*
492 *Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2011,*
493 *Athens, Greece, September 5-9, 2011. Proceedings, Part I 11*, pp. 12–27. Springer, 2011.
- 494
495 Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal
496 Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from
497 human preferences. In *International Conference on Artificial Intelligence and Statistics*, pp.
4447–4455. PMLR, 2024.
- 498
499 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
500 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 501
502 Róbert Busa-Fekete, Balázs Szörényi, Paul Weng, Weiwei Cheng, and Eyke Hüllermeier. Preference-
503 based reinforcement learning: evolutionary direct policy search using a preference-based racing
504 algorithm. *Machine learning*, 97:327–351, 2014.
- 505
506 Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier
507 Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems
508 and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint*
arXiv:2307.15217, 2023.
- 509
510 Xiaoyu Chen, Han Zhong, Zhuoran Yang, Zhaoran Wang, and Liwei Wang. Human-in-the-loop:
511 Provably efficient preference-based reinforcement learning with general function approximation.
512 In *International Conference on Machine Learning*, pp. 3773–3793. PMLR, 2022.
- 513
514 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
515 reinforcement learning from human preferences. *Advances in neural information processing*
516 *systems*, 30, 2017.
- 517
518 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model
519 alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- 520
521 Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre
522 Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online
523 ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- 524
525 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with
526 deep energy-based policies. In *ICML*, 2017.
- 527
528 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
529 maximum entropy deep reinforcement learning with a stochastic actor. In *International conference*
530 *on machine learning*, pp. 1861–1870. Pmlr, 2018.
- 531
532 Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without
533 reference model. *arXiv preprint arXiv:2403.07691*, 2024.
- 534
535 B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani,
536 and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE*
537 *transactions on intelligent transportation systems*, 23(6):4909–4926, 2021.
- 538
539 W Bradley Knox and Peter Stone. Tamer: Training an agent manually via evaluative reinforcement.
In *2008 7th IEEE international conference on development and learning*, pp. 292–297. IEEE,
2008.
- 538
539 Pushmeet Kohli, Mahyar Salek, and Greg Stoddard. A fast bandit algorithm for recommendation to
users with heterogenous tastes. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
volume 27, pp. 1135–1141, 2013.

- 540 Gengxu Li, Tingyu Xia, Yi Chang, and Yuan Wu. Length-controlled margin-based preference
541 optimization without reference model. *arXiv preprint arXiv:2502.14643*, 2025.
- 542
- 543 Runze Liu, Fengshuo Bai, Yali Du, and Yaodong Yang. Meta-reward-net: Implicitly differentiable
544 reward learning for preference-based reinforcement learning. *Advances in Neural Information
545 Processing Systems*, 35:22270–22284, 2022.
- 546
- 547 Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-
548 free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235, 2024.
- 549
- 550 Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1,
551 pp. 2, 2000.
- 552
- 553 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
554 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
555 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
27744, 2022.
- 556
- 557 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality in
558 direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024.
- 559
- 560 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
561 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances
562 in Neural Information Processing Systems*, 36:53728–53741, 2023.
- 563
- 564 Rafael Rafailov, Joey Hejna, Ryan Park, and Chelsea Finn. From R to Q*: Your language model is
565 secretly a Q-function. *arXiv preprint arXiv:2404.12358*, 2024.
- 566
- 567 Aadirupa Saha, Aldo Pacchiano, and Jonathan Lee. Dueling rl: Reinforcement learning with trajectory
568 preferences. In *International conference on artificial intelligence and statistics*, pp. 6263–6289.
569 PMLR, 2023.
- 570
- 571 Marc Schoenauer, Riad Akrou, Michele Sebag, and Jean-Christophe Souplet. Programming by
572 feedback. In *International Conference on Machine Learning*, pp. 1503–1511. PMLR, 2014.
- 573
- 574 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
575 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 576
- 577 Arsalan Sharifnassab, Saber Salehkaleybar, Sina Ghiassian, Surya Kanoria, and Dale Schuurmans.
578 Soft preference optimization: Aligning language models to expert distributions. *arXiv preprint
579 arXiv:2405.00747*, 2024.
- 580
- 581 Daniel A Spielman. Algorithms, graph theory, and linear equations in laplacian matrices. In
582 *Proceedings of the International Congress of Mathematicians 2010 (ICM 2010) (In 4 Volumes) Vol.
583 I: Plenary Lectures and Ceremonies Vols. II–IV: Invited Lectures*, pp. 2698–2722. World Scientific,
584 2010.
- 585
- 586 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
587 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in
588 neural information processing systems*, 33:3008–3021, 2020.
- 589
- 590 Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaxi-
591 malist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*,
592 2024.
- 593
- 594 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
595 In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033.
596 IEEE, 2012.
- 597
- 598 Garrett Warnell, Nicholas Waytowich, Vernon Lawhern, and Peter Stone. Deep tamer: Interactive
599 agent shaping in high-dimensional state spaces. In *Proceedings of the AAAI conference on artificial
600 intelligence*, volume 32, 2018.

- 594 Aaron Wilson, Alan Fern, and Prasad Tadepalli. A bayesian approach for policy learning from
595 trajectory preference queries. *Advances in neural information processing systems*, 25, 2012.
596
- 597 Runzhe Wu and Wen Sun. Making rl with preference-based feedback efficient via randomization.
598 *arXiv preprint arXiv:2310.14554*, 2023.
- 599 Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and
600 Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q^* -approximation
601 for sample-efficient rlhf. *arXiv preprint arXiv:2405.21046*, 2024.
602
- 603 Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton
604 Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm
605 performance in machine translation. *arXiv preprint arXiv:2401.08417*, 2024.
- 606 Yichong Xu, Ruosong Wang, Lin Yang, Aarti Singh, and Artur Dubrawski. Preference-based
607 reinforcement learning with finite-time guarantees. *Advances in Neural Information Processing
608 Systems*, 33:18784–18794, 2020.
- 609 Chunqiu Zeng, Qing Wang, Shekoofeh Mokhtari, and Tao Li. Online context-aware recommendation
610 with time varying multi-armed bandit. In *Proceedings of the 22nd ACM SIGKDD international
611 conference on Knowledge discovery and data mining*, pp. 2025–2034, 2016.
612
- 613 Wenhao Zhan, Masatoshi Uehara, Wen Sun, and Jason D Lee. Provable reward-agnostic preference-
614 based reinforcement learning. *arXiv preprint arXiv:2305.18505*, 2023.
- 615 Qining Zhang and Lei Ying. Zeroth-order policy gradient for reinforcement learning from human
616 feedback without reward inference. *arXiv preprint arXiv:2409.17401*, 2024.
617
- 618 Qining Zhang, Honghao Wei, and Lei Ying. Reinforcement learning from human feedback with-
619 out reward inference: Model-free algorithm and instance-dependent analysis. *arXiv preprint
620 arXiv:2406.07455*, 2024.
- 621 Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human
622 feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*,
623 pp. 43037–43067. PMLR, 2023.
624
- 625 Banghua Zhu, Michael I Jordan, and Jiantao Jiao. Iterative data smoothing: Mitigating reward
626 overfitting and overoptimization in rlhf. *arXiv preprint arXiv:2401.16335*, 2024.
- 627 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
628 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv
629 preprint arXiv:1909.08593*, 2019.
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A PROOF OF THE THEOREM 5.2

Assumption A.1 (Bradley–Terry preferences on the *soft*-optimal Q -function). Let $Q^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ be the soft-optimal state-action value function of the MDP, i.e.,

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \mathcal{H}(\pi(\cdot | s_t))) \mid s_0 = s, a_0 = a \right],$$

where $\mathcal{H}(\cdot)$ is the entropy function and λ is the entropy coefficient. Assume that there exists a parameter $\beta > 0$ such that, for every $s \in \mathcal{S}$ and any $a, b \in \mathcal{A}$,

$$P^*(a \succ b | s) = \sigma(\beta [Q^*(s, a) - Q^*(s, b)]),$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}.$$

Theorem A.2 (Preference loss recovers the optimal policy). Fix a state $s \in \mathcal{S}$ and abbreviate $Q_a^* := Q^*(s, a)$. Suppose that Theorem A.1 holds. Under uniform sampling of ordered pairs $(a, b) \sim \text{Unif}(\mathcal{A}^2)$ and the tabular full-support parameterization $\ell_a = \log \pi(a | s)$ ($\sum_a e^{\ell_a} = 1$, $e^{\ell_a} > 0$), consider the preference loss

$$\mathcal{L}_{\text{pref}}(\ell) = - \frac{1}{|\mathcal{A}|^2} \sum_{(a,b) \in \mathcal{A}^2} P^*(a \succ b | s) \log \sigma(\alpha(\ell_a - \ell_b)),$$

$$\alpha > 0. \quad (7)$$

This loss is strictly convex on the set of full-support policies and is minimized uniquely at

$$\pi_*(a | s) = \frac{\exp(\frac{\beta}{\alpha} Q_a^*)}{\sum_{a' \in \mathcal{A}} \exp(\frac{\beta}{\alpha} Q_{a'}^*)}. \quad (8)$$

Furthermore, this Gibbs distribution coincides with the global maximizer of the entropy-regularized RL (or SAC objective) when $\lambda = \alpha/\beta$:

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \lambda \mathcal{H}(\pi(\cdot | s_t))) \right]. \quad (9)$$

Proof. Because the policy is tabular, we fix the state s , and introduce the log-policy vector $\ell = (\ell_a)_{a \in \mathcal{A}}$. Define the policy and the Bradley–Terry probabilities as follows:

$$P_{ab}(\ell) := \sigma(\alpha(\ell_a - \ell_b)), \quad P_{ab}^* := \sigma(\beta(Q_a^* - Q_b^*)), \quad a, b \in \mathcal{A}.$$

First, we reformulate the loss of the theorem. For this purpose, we consider two cases:

1. When $a = b$: In this case the two logits coincide, so $P_{aa}^* = P_{aa} = \sigma(0) = \frac{1}{2}$; hence in this case each summand equals $-\frac{1}{2} \log \frac{1}{2} = \frac{\log 2}{2}$. Summing over the $|\mathcal{A}|$ therefore contributes the constant $\frac{|\mathcal{A}| \log 2}{2|\mathcal{A}|^2}$ in the loss.
2. For any two different actions $a \neq b$ the ordered pairs (a, b) and (b, a) both appear. Because $\sigma(z) + \sigma(-z) = 1$, we have the identities $P_{ba} = 1 - P_{ab}$ and $P_{ba}^* = 1 - P_{ab}^*$. Grouping those two ordered terms gives the compact expression

$$\mathcal{L}(\ell) := - \frac{1}{|\mathcal{A}|^2} \sum_{\{a,b\} \in \mathcal{A}, a \neq b} \left[P_{ab}^* \log P_{ab}(\ell) + P_{ba}^* \log P_{ba}(\ell) \right] \quad (10)$$

Therefore, $\mathcal{L}_{\text{pref}} = \frac{|\mathcal{A}| \log 2}{2|\mathcal{A}|^2} + \mathcal{L}$. Since the additive constant is not used in the optimization, it can be discarded. Therefore, we may optimize \mathcal{L} instead of $\mathcal{L}_{\text{pref}}$.

Now we characterize the stationary points. For this purpose, we compute the partial derivative of \mathcal{L} with respect to the ℓ_k , $\frac{\partial \mathcal{L}}{\partial \ell_k}$. Based on Lemma A.3 only the terms that contain k depend on ℓ_k , so

$$\frac{\partial \mathcal{L}}{\partial \ell_k} = -\frac{\alpha}{|\mathcal{A}|^2} \sum_{b \neq k} [P_{kb}^* - P_{kb}(\ell)]. \quad (11)$$

A stationary point satisfies $\sum_{b \neq k} (P_{kb} - P_{kb}^*) = 0$ for every k .

Set

$$\ell_a = c + \frac{\beta}{\alpha} Q_a^*, \quad \forall a \in \mathcal{A} \quad (12)$$

Then for every a, b ,

$$P_{ab}(\ell) = \sigma(\alpha(\ell_a - \ell_b)) = \sigma(\beta(Q_a^* - Q_b^*)) = P_{ab}^*,$$

and equation 11 is equal to zero with $\ell_a = c + \frac{\beta}{\alpha} Q_a^*$.

Now using $\sum_a e^{\ell_a} = 1$ with equation 12 gives

$$e^c = \left(\sum_a \exp\left(\frac{\beta}{\alpha} Q_a^*\right) \right)^{-1}.$$

Hence, the stationary point is as follows:

$$\pi_*(a | s) = \frac{\exp\left(\frac{\beta}{\alpha} Q_a^*\right)}{\sum_{a' \in \mathcal{A}} \exp\left(\frac{\beta}{\alpha} Q_{a'}^*\right)}. \quad (13)$$

If we write the KKT conditions of the loss and derive the value of the Lagrange multiplier λ , λ will be zero. Hence, the above stationary point is valid. See Lemma A.4 for the details.

Computing Hessian: To compute the Hessian we define w_{ab} as follows:

$$w_{ab} := P_{ab}(\ell)P_{ba}(\ell) = P_{ab}(\ell)[1 - P_{ab}(\ell)] > 0.$$

Using

$$\frac{\partial P_{ab}}{\partial \ell_a} = +\alpha w_{ab}, \quad \frac{\partial P_{ab}}{\partial \ell_b} = -\alpha w_{ab}.$$

Now, if we differentiate equation 11 once more. For ($i \neq j$):

$$\frac{\partial^2 \mathcal{L}_p}{\partial \ell_j \partial \ell_i} = -\frac{\alpha}{|\mathcal{A}|^2} [\alpha w_{ij}] = -\frac{\alpha^2}{|\mathcal{A}|^2} w_{ij}.$$

For ($i = j$):

$$\frac{\partial^2 \mathcal{L}_p}{\partial \ell_i^2} = -\frac{\alpha}{|\mathcal{A}|^2} \sum_{b \neq i} (-\alpha w_{ib}) = \frac{\alpha^2}{|\mathcal{A}|^2} \sum_{b \neq i} w_{ib}.$$

Now using the above derivations, we write the matrix form of Hessian.

$$\mathbf{H} = \frac{\alpha^2}{|\mathcal{A}|^2} (\mathbf{D} - \mathbf{W}) \quad \begin{aligned} W_{ij} &= w_{ij} \ (i \neq j), \quad W_{ii} = 0, \\ D_{ii} &= \sum_{b \neq i} w_{ib}. \end{aligned}$$

To prove that \mathcal{L} is strictly convex, we should prove that \mathbf{H} (or \mathbf{L}) is positive-definite. The matrix $\mathbf{L} = \mathbf{D} - \mathbf{W}$ is a weighted graph Laplacian of the complete graph on \mathcal{A} as its off-diagonal entries are negative, diagonals are positive, and each row sums to zero (Spielman, 2010).

For a matrix \mathbf{L} to be positive definite, we should have for any $v \in \mathbb{R}^{|\mathcal{A}|}$, $v^\top \mathbf{L}v > 0$. In our case, one has

$$v^\top \mathbf{L}v = \frac{1}{2} \sum_{i,j} w_{ij} (v_i - v_j)^2. \quad (14)$$

Identity equation 14 follows from expanding $v^\top (\mathbf{D} - \mathbf{W})v$ and re-grouping terms (see Spielman (2010) for more details). Because every weight $w_{ij} > 0$, the RHS is non-negative, hence it is always equal to or bigger than zero. Therefore \mathbf{L} is positive-semidefinite ($\mathbf{L} \succeq 0$) and equals to zero iff

$$v_1 = \dots = v_{|\mathcal{A}|}.$$

In other words, only subspace $\text{span}\{\mathbf{1}\} = \{a\mathbf{1} : a \in \mathbb{R}, a \neq 0\}$, whose members have all coordinates equal ($v_1 = \dots = v_{|\mathcal{A}|}$) makes $v^\top \mathbf{L}v$ equal to zero. Now, we prove that given the constraint imposed on our problem, $v^\top \mathbf{L}v$ cannot be equal to zero.

In general unconstrained optimization, v in $v^\top \mathbf{L}v$, shows all possible directions in $\mathbb{R}^{|\mathcal{A}|}$. In our case, the optimization is constrained, and the function \mathcal{L} is restricted to a constraint set C (the probability simplex: $\sum_{a \in \mathcal{A}} e^{l_a} - 1 = 0$), hence the condition $v^\top \mathbf{H}v \geq 0$ is only required for vectors v in the tangent space of C at ℓ . This is because the tangent space of a convex set C at any $v \in C$ is the set of feasible directions within C (Absil et al., 2009).

The parameter set of the \mathbf{H} (or accordingly \mathbf{L}) is

$$\mathcal{E} := \left\{ \ell \in \mathbb{R}^{|\mathcal{A}|} : \sum_a e^{l_a} = 1, e^{l_a} > 0 \right\}$$

We define $g(\ell) = \sum_{a \in \mathcal{A}} e^{l_a} - 1 = 0$ and its gradient $\nabla g(\ell) = e^\ell := (e^{l_1}, \dots, e^{l_{|\mathcal{A}|}})^\top > 0$. A displacement $v \in \mathbb{R}^{|\mathcal{A}|}$ is *feasible* iff it is in the tangent space (denote it with $T_\ell \mathcal{E}$) that is $\nabla g(\ell) \cdot v = 0$. Hence,

$$e^\ell \cdot v = 0.$$

Now consider a vector in $\text{span}\{\mathbf{1}\}$. For any $v = a\mathbf{1}$ with $a \neq 0$,

$$e^\ell \cdot v = a e^\ell \cdot \mathbf{1} = a \sum_{b \in \mathcal{A}} e^{l_b} = a \neq 0.$$

Hence, $v \notin T_\ell \mathcal{E}$. Hence, for every *feasible* $v \neq 0$,

$$v^\top \mathbf{L}v > 0 \implies v^\top \mathbf{H}v = \frac{\alpha^2}{|\mathcal{A}|^2} v^\top \mathbf{L}v > 0.$$

Thus, the Hessian is positive-definite along all feasible directions, which establishes the strict convexity of the preference loss on the full support tabular policy.

Another way to proof the uniqueness of the solution: Assume, for contradiction, that there exists another log-policy vector $\tilde{\ell} \in \mathcal{E}$ that also satisfies the stationarity system $\sum_{b \neq k} (P_{kb} - P_{kb}^*) = 0$, $\forall k$ and the normalization $\sum_a e^{\tilde{l}_a} = 1$. For every ordered pair (k, j) the argument leading to equation ?? then gives $P_{kj}(\tilde{\ell}) = P_{kj}^*$ as well. Because σ is strictly increasing, this implies

$$\tilde{l}_k - \tilde{l}_j = l_k - l_j, \quad \forall k, j \in \mathcal{A},$$

hence $\tilde{\ell} = \ell + \delta \mathbf{1}$ for some $\delta \in \mathbb{R} \setminus \{0\}$. But then

$$\sum_a e^{\tilde{l}_a} = e^\delta \sum_a e^{l_a} = e^\delta \neq 1,$$

contradicting the constraint that every feasible ℓ must satisfy $\sum_a e^{l_a} = 1$. Therefore $\delta = 0$ and $\tilde{\ell} = \ell$, proving that the stationary point is unique.

Connection with soft actor-critic: For the same fixed state consider

$$J_\tau(\pi) := \mathbb{E}_{a \sim \pi}[Q_a^*] + \tau \mathcal{H}(\pi), \quad \mathcal{H}(\pi) := - \sum_a \pi(a) \log \pi(a).$$

Introducing a Lagrange multiplier λ for $\sum_a \pi(a) = 1$ gives $Q_a^* - \tau(\log \pi(a) + 1) - \lambda = 0$, hence $\pi(a) \propto \exp(Q_a^*/\tau)$. Normalization produces

$$\pi_{\text{SAC}}(a | s) = \frac{\exp(Q_a^*/\tau)}{\sum_{a'} \exp(Q_{a'}^*/\tau)}.$$

Choosing $\tau = \alpha/\beta$ recovers equation 13, so the minimizer of \mathcal{L} coincides with the soft actor-critic solution with temperature $\tau = \alpha/\beta$.

□

Lemma A.3 (Gradient of the unordered-pair preference loss). *Let*

$$\mathcal{L}(\ell) := -\frac{1}{|\mathcal{A}|^2} \sum_{\{a,b\} \in \mathcal{A}, a \neq b} \left[P_{ab}^* \log P_{ab}(\ell) + P_{ba}^* \log P_{ba}(\ell) \right],$$

where P_{ab}^* , $P_{ab}(\ell)$ and ℓ are defined in Theorem A.2. Then, for every action $k \in \mathcal{A}$,

$$\frac{\partial \mathcal{L}}{\partial \ell_k} = -\frac{\alpha}{|\mathcal{A}|^2} \sum_{b \neq k} [P_{kb}^* - P_{kb}(\ell)]. \quad (15)$$

Proof. For each unordered pair $\{a, b\}$, define

$$g_{ab}(\ell) := P_{ab}^* \log P_{ab}(\ell) + P_{ba}^* \log P_{ba}(\ell).$$

Because $\frac{d}{dz} \log \sigma(z) = 1 - \sigma(z)$ and $\partial(\ell_a - \ell_b)/\partial \ell_k = \mathbf{1}\{k = a\} - \mathbf{1}\{k = b\}$,

$$\frac{\partial}{\partial \ell_k} \log P_{ab} = \alpha(1 - P_{ab})[\mathbf{1}\{k = a\} - \mathbf{1}\{k = b\}],$$

$$\frac{\partial}{\partial \ell_k} \log P_{ba} = \alpha P_{ab}[\mathbf{1}\{k = b\} - \mathbf{1}\{k = a\}].$$

Since $P_{ba}^* = 1 - P_{ab}^*$ and $P_{ba} = 1 - P_{ab}$,

$$\frac{\partial g_{ab}}{\partial \ell_k} = \alpha(\mathbf{1}\{k = a\} - \mathbf{1}\{k = b\})[P_{ab}^* - P_{ab}]. \quad (16)$$

Insert equation 16 into the loss and sum over all unordered pairs that contain k :

$$\frac{\partial \mathcal{L}}{\partial \ell_k} = -\frac{\alpha}{|\mathcal{A}|^2} \sum_{\{a,b\}} (\mathbf{1}\{k = a\} - \mathbf{1}\{k = b\}) [P_{ab}^* - P_{ab}].$$

If $k = a$ (and $b > k$) the indicator equals $+1$; if $k = b$ (with $a < k$) it equals -1 . Using again the symmetry $P_{ak}^* = 1 - P_{ka}^*$ and $P_{ak} = 1 - P_{ka}$, the negative sign flips the difference so that both cases contribute the same quantity $P_{kb}^* - P_{kb}$. Hence

$$\frac{\partial \mathcal{L}}{\partial \ell_k} = -\frac{\alpha}{|\mathcal{A}|^2} \sum_{b \neq k} [P_{kb}^* - P_{kb}(\ell)],$$

completing the proof. □

Lemma A.4 (The KKT multiplier). *Consider the constrained minimization of the unordered-pair preference loss*

$$\min_{\ell \in \mathbb{R}^{|\mathcal{A}|}} \mathcal{L}(\ell) \quad \text{s.t.} \quad g(\ell) := \sum_{a \in \mathcal{A}} e^{\ell_a} - 1 = 0,$$

with \mathcal{L} defined in equation 10. Let $\lambda \in \mathbb{R}$ be the Lagrange multiplier associated with the normalization constraint. At every KKT point (ℓ, λ) one necessarily has

$$\lambda = 0.$$

864 *Proof.* The KKT stationarity condition for each action $k \in \mathcal{A}$ is

$$865 \quad -\frac{\alpha}{|\mathcal{A}|^2} \sum_{b \neq k} [P_{kb}^* - P_{kb}(\ell)] + \lambda e^{\ell_k} = 0. \quad (16)$$

867 Summing equation 16 over all k gives

$$868 \quad -\frac{\alpha}{|\mathcal{A}|^2} \sum_k \sum_{b \neq k} [P_{kb}^* - P_{kb}(\ell)] + \lambda \sum_k e^{\ell_k} = 0. \quad (17)$$

869 Because the log-policy variables satisfy the equality constraint $\sum_k e^{\ell_k} = 1$, the second term in equation 17 sums to λ .

870 Rewrite the double sum by grouping every *ordered* pair (k, b) with its reverse (b, k) :

$$871 \quad \sum_k \sum_{b \neq k} [P_{kb}^* - P_{kb}] = \sum_{\substack{k, b \in \mathcal{A} \\ k < b}} [(P_{kb}^* - P_{kb}) + (P_{bk}^* - P_{bk})].$$

872 Using the fact that $P_{kb}^* + P_{bk}^* = 1$ and $P_{kb} + P_{bk} = 1$, each term in the bracket equals to $1 - 1 = 0$. Therefore, the entire double sum is zero, and we can imply that $\lambda = 0$.

873 \square

874 B TRAJECTORY-LEVEL ANALYSIS OF DFA

875 Let \mathcal{T}_H be the (finite) set of all length- H trajectories $\tau = (s_0, a_0, \dots, s_{H-1}, a_{H-1})$ that can be generated by the MDP.

876 **Assumption B.1** (Trajectory-level Bradley-Terry model). Let

$$877 \quad G^*(\tau) := \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t)$$

878 be the *return* of trajectory τ . There exists $\beta > 0$ such that for every pair $\tau_1, \tau_2 \in \mathcal{T}_H$

$$879 \quad P^*(\tau_1 \succ \tau_2) = \sigma(\beta [G^*(\tau_1) - G^*(\tau_2)]), \quad \sigma(z) = \frac{1}{1+e^{-z}}.$$

880 B.1 TRAJECTORY PREFERENCE LOSS

881 Parameterise a *trajectory-tabular* policy by one log-likelihood per path, $L_\tau = \log \pi_\theta(\tau)$, subject to the simplex constraint $\sum_{\tau \in \mathcal{T}_H} e^{L_\tau} = 1$, $e^{L_\tau} > 0$. For ordered trajectory pairs sampled uniformly from \mathcal{T}_H^2 define the loss

$$882 \quad \mathcal{L}_{\text{traj}}(L) := -\frac{1}{|\mathcal{T}_H|^2} \sum_{\tau_1, \tau_2 \in \mathcal{T}_H} P^*(\tau_1 \succ \tau_2) \log \sigma(\alpha [L_{\tau_1} - L_{\tau_2}]), \quad (18)$$

883 with parameter $\alpha > 0$.

884 **Theorem B.2** (Optimal policy for trajectory loss). *Assume Theorem B.1 and uniform sampling of ordered trajectory pairs. The loss equation 18 is strictly convex on the probability simplex $\{L : \sum_\tau e^{L_\tau} = 1\}$ and attains its unique minimum at the Gibbs distribution*

$$885 \quad \pi_*(\tau) = \frac{\exp(\frac{\beta}{\alpha} G^*(\tau))}{\sum_{\tau' \in \mathcal{T}_H} \exp(\frac{\beta}{\alpha} G^*(\tau'))}. \quad (19)$$

886 *The proof is similar to the proof of Theorem 5.2*

B.2 CONNECTION BETWEEN THE STATE-WISE AND TRAJECTORY-WISE OPTIMA

We adopt the standard soft Bellman optimality equations (See Haarnoja et al. (2017) for the proofs):

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)}[V^*(s')], \quad (20)$$

$$V^*(s) = \alpha/\beta \log \sum_{a' \in \mathcal{A}} \exp\left(\frac{\beta}{\alpha} Q^*(s, a')\right). \quad (21)$$

Let $\alpha, \beta > 0$ be the preference and BT-scale parameters. The optimal (state-wise) policy can be written equivalently as

$$\pi_{\text{st}}(a | s) = \frac{\exp\left(\frac{\beta}{\alpha} Q^*(s, a)\right)}{\sum_{a'} \exp\left(\frac{\beta}{\alpha} Q^*(s, a')\right)} = \exp\left(\frac{\beta}{\alpha} (Q^*(s, a) - V^*(s))\right). \quad (22)$$

Consider a finite-horizon trajectory $\tau = (s_0, a_0, \dots, s_{H-1}, a_{H-1})$ generated by the deterministic dynamics $s_{t+1} = f(s_t, a_t)$, with fixed initial state s_0 and terminal boundary condition $V^*(s_H) = 0$. Define the trajectory probability induced by the state-wise policy,

$$\pi_{\text{traj}}(\tau) := \prod_{t=0}^{H-1} \pi_{\text{st}}(a_t | s_t). \quad (23)$$

Taking logs and using equation 22,

$$\log \pi_{\text{traj}}(\tau) = \frac{\beta}{\alpha} \sum_{t=0}^{H-1} (Q^*(s_t, a_t) - V^*(s_t)). \quad (24)$$

Under the assumptions $\gamma = 1$ and deterministic transitions, the soft Bellman equation equation 20 reduces to

$$Q^*(s_t, a_t) = r_t + V^*(s_{t+1}), \quad (25)$$

hence

$$Q^*(s_t, a_t) - V^*(s_t) = r_t + V^*(s_{t+1}) - V^*(s_t). \quad (26)$$

Summing over $t = 0, \dots, H - 1$ telescopes the value terms:

$$\sum_{t=0}^{H-1} (Q^*(s_t, a_t) - V^*(s_t)) = \sum_{t=0}^{H-1} r_t + V^*(s_H) - V^*(s_0) = \sum_{t=0}^{H-1} r_t - V^*(s_0). \quad (27)$$

Combining equation 24 and equation 27 gives

$$\log \pi_{\text{traj}}(\tau) = \frac{\beta}{\alpha} \left(\sum_{t=0}^{H-1} r_t - V^*(s_0) \right). \quad (28)$$

Therefore, under these conditions, for any trajectory τ starting from a fixed s_0 we have

$$\pi_{\text{traj}}(\tau) = \exp\left(\frac{\beta}{\alpha} (G(\tau) - V^*(s_0))\right) = C(s_0) \exp\left(\frac{\beta}{\alpha} G(\tau)\right) \propto \exp\left(\frac{\beta}{\alpha} G(\tau)\right),$$

where $G(\tau) = \sum_{t=0}^{H-1} r_t$ and $C(s_0) = \exp(-\frac{\beta}{\alpha} V^*(s_0))$ depends only on the common start state.

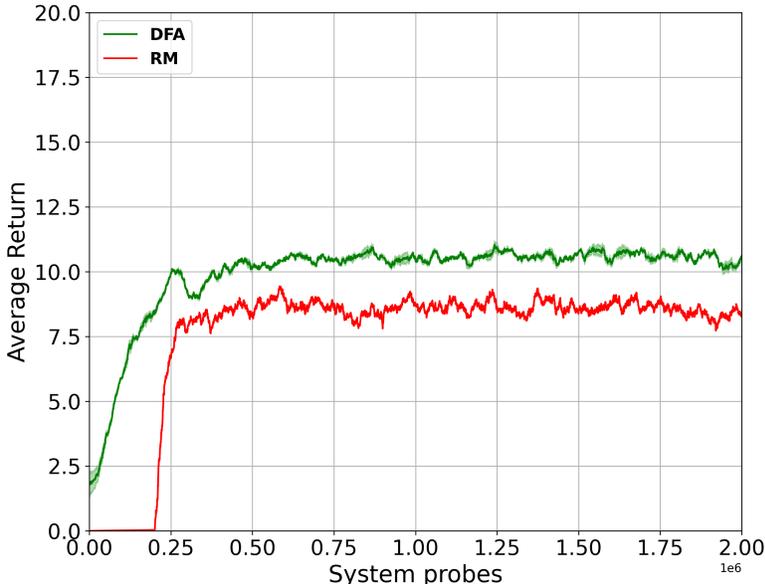
This establishes that the trajectory probabilities induced by the state-wise soft-optimal policy are proportional to $\exp((\beta/\alpha)G(\tau))$ (for fixed s_0), completing the connection between the state-wise and trajectory-wise formulations.

Moreover, for any two trajectories τ^+, τ^- with the same s_0 ,

$$\alpha [\log \pi_{\text{traj}}(\tau^+) - \log \pi_{\text{traj}}(\tau^-)] = \beta [G(\tau^+) - G(\tau^-)],$$

so the DFA trajectory-level preference computed from π_{traj} reduces exactly to the Bradley-Terry model over undiscounted returns:

$$P_\theta(\tau^+ \succ \tau^-) = \sigma\left(\beta [G(\tau^+) - G(\tau^-)]\right).$$

Figure 4: GridWorld results with size 10*10. ($M = 500$)

C EXPERIMENTS

Walker is a planar biped with four actuated joints that must walk without tipping over; Hopper is a one-leg, three-joint robot that learns to hop forward; Humanoid is a 17-joint 3D figure that must walk quickly while remaining upright; Swimmer is a three-link snake that propels itself through a viscous medium; Inverted Pendulum tasks a cart with balancing an upright pole; and MountainCar Continuous challenges a car trapped between two hills to climb the right hill by building momentum.

For GridWorld game, the grid is 5*5, where the agent starts at position 2*2. All methods use a tabular softmax policy parameterization, where each state-action pair has a corresponding logit parameter. For the reward model in RM+PPO, we use a simple tabular representation that assigns a value to each state-action pair. All methods are trained for the same number of environment interactions to ensure fair comparison. We use Adam optimizer with learning rate 3×10^{-2} across all methods.

Below we illustrate the results for more complex environments and different numbers of M mentioned in Section 6. We run the algorithms for 5 different seeds {3, 1, 14, 4, 50}. We considered the default gym horizon for all environments.

Figure 7 shows the performance of DFA on Humanoid environments compared with OraclePPO and RM methods. For the panel query, with used $M=10$. For the RM method, we spent 2000 system probes for the reward model training. In tasks whose rewards are easy to model, the performance gap between RM and Oracle PPO is typically small; we therefore chose Humanoid, where the reward is harder to model than in the other environments.

We utilized a Linux server with Intel Xeon CPU E5-2680 v3 (24 cores) operating at 2.50GHz with 377 GB DDR4 of memory and Nvidia Titan X Pascal GPU. The computation was distributed over 48 threads to ensure a relatively efficient run time. In our control-task experiments, DFA required more wall-clock time than SAC. For example, in the Pendulum, running 10 million system probes took 8 hours (on average) with DFA compared to 6.5 hours with SAC. In the Swimmer environment, SAC completed in 8 hours, while DFA took 9 hours. Although DFA generally requires more wall-clock time per step, in some environments (e.g., MountainCar, Swimmer) it converges in fewer steps. As a result, the increased per-step runtime does not significantly impact its overall efficiency.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047

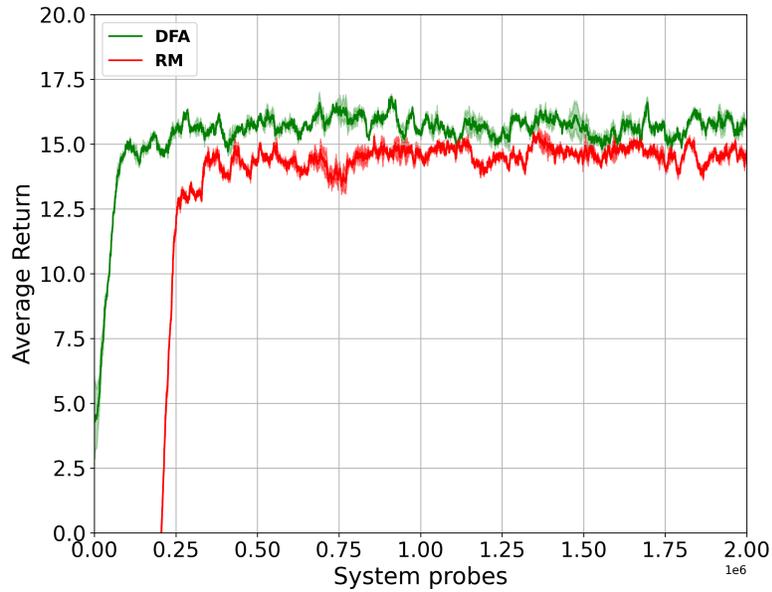


Figure 5: GridWorld results with size $20 \times 20 (M = 100)$.

1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072

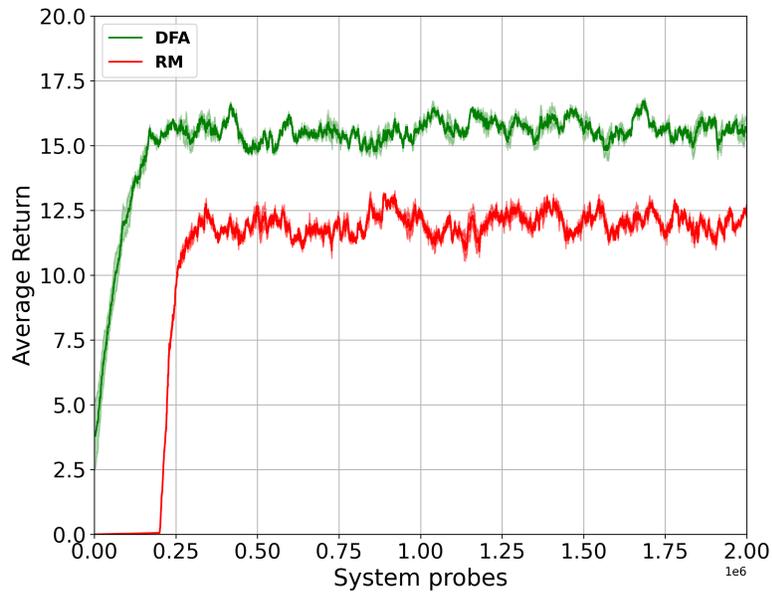


Figure 6: GridWorld results with size $20 \times 20 (M = 1)$.

1073
1074
1075
1076
1077
1078
1079

For hyperparameter tuning, we performed a grid search, systematically exploring a predefined range of values for each parameter. In the following tables, we provide the fine-tuned parameters for each algorithm and method. Batch sizes are considered the same for all algorithms. The discount factor is also set to 0.99 for all the runs.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Table 1: Default hyper-parameters for all algorithms used in the 5×5 GridWorld experiments

Algorithm	Hyper-parameter	Value
ZPG	T (iterations)	1000000
	N (pairs / iter)	1 - 10
	M (votes / query)	1000
	μ (perturbation radius)	0.1
	α (learning-rate)	0.05
	trim (prob. clip)	10^{-2}
RM-PPO	traj_pairs (pretaining)	5000
	ppo_iters	1000
	β_{KL}	0.1
	γ (discount)	1.0
	λ (GAE)	0.95
DFA (on-policy)	α (temperature)	$1 \times 10^{-3} - 1 \times 10^{-6}$
	$N_{\text{pairs/iter}}$	1
	iters	100000
Oracle-PPO	ppo_iters	100000
	β_{KL}	0.1
	γ (discount)	1.0
	λ (GAE)	0.95

Table 2: Hyper-parameters for SAC and DFA across all evaluated environments

Alg.	Hyper-parameter	Walker2d	Hopper	Swimmer	Humanoid	MountainCarC	Pendulum
SAC	Hidden layer size	64	64	64	64	64	64
	Policy learning-rate	1×10^{-3}					
	Q learning-rate	1×10^{-3}					
	Batch size	256	256	256	256	256	256
	Replay-buffer capacity	20 000	20 000	20 000	20 000	20 000	20 000
	Entropy temperature λ	0.1	0.2	0.01	0.01	0.1	0.2
	Discount factor γ	0.99	0.99	0.99	0.99	0.99	0.99
	Soft-update coefficient τ	0.1	0.005	0.1	0.1	0.01	0.005
	# parallel envs N_{env}	32	32	32	32	32	32
	Training episodes	50 000	50 000	50 000	50 000	50 000	50 000
	DFA	Hidden layer size	64	64	64	64	64
Policy learning-rate		1×10^{-3}					
Q learning-rate		1×10^{-3}					
Batch size		256	256	256	256	256	256
Replay-buffer capacity		20 000	20 000	20 000	20 000	20 000	20 000
Entropy temperature λ		0.01	0.1	0.01	0.01	0.01	0.01
Temperature α		0.2	0.2	0.3	0.2	0.4	0.2
Discount factor γ		0.99	0.99	0.99	0.99	0.99	0.99
Soft-update coefficient τ		0.1	0.005	0.1	0.1	0.01	0.005
# parallel envs N_{env}		32	32	32	32	32	32
Training episodes		50 000	50 000	50 000	50 000	50 000	50 000

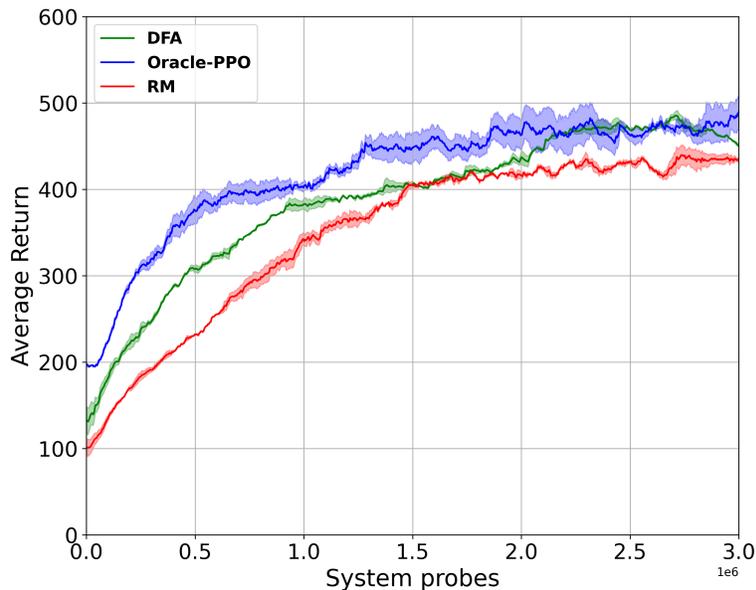


Figure 7: Humanoid results with horizon 1000.

D BROADER IMPACTS

DFA aims to make reinforcement learning from human feedback more sample-efficient by blending numeric rewards with pairwise preferences. Positive impacts include lowering annotation costs, enabling faster prototyping of assistive robots, and providing a simple baseline for preference-centric alignment research. However, the method also amplifies whatever biases or inconsistencies are present in the collected preferences: if early Q estimates or human labels encode unfair or unsafe behavior, DFA may reinforce those patterns more quickly than reward-only training. Because DFA can learn from very small amounts of feedback, malicious or accidental injection of adversarial comparisons could steer policies toward harmful objectives—especially in safety-critical domains such as autonomous driving or content recommendation. The work uses only simulated environments and involves no personal data; nevertheless, broader deployment should respect fairness guidelines and, when real users provide feedback, comply with relevant privacy regulations.

E USE OF LARGE LANGUAGE MODELS

We used Large Language Models (LLMs) to aid or polish the manuscript text. Specifically, LLMs were used to improve grammar, phrasing, and clarity of exposition; they were also used for code debugging.

F REBUTTAL EXPERIMENT

We added one more experiment during the rebuttal discussion on a new environment in MetaWorld, requested by the reviewers. The experiment setup is consistent with the paper Liu et al. (2022). We will update this section for the camera-ready version with more experiments.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

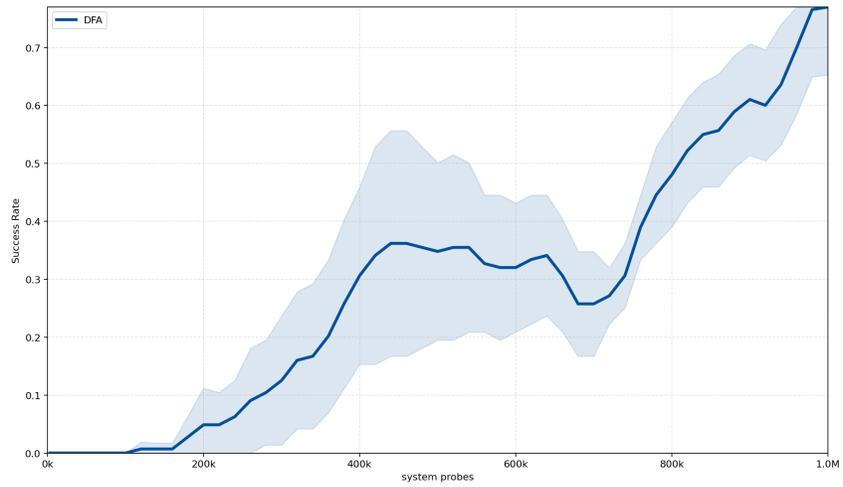


Figure 8: Rebuttal Hammer results.