

---

# Reward Model Learning vs. Direct Policy Optimization: A Comparative Analysis of Learning from Human Preferences

---

Andi Nika<sup>1</sup> Debmalya Mandal<sup>2</sup> Parameswaran Kamalaruban<sup>3</sup> Georgios Tzannetos<sup>1</sup>  
Goran Radanović<sup>1</sup> Adish Singla<sup>1</sup>

## Abstract

In this paper, we take a step towards a deeper understanding of learning from human preferences by systematically comparing the paradigm of reinforcement learning from human feedback (RLHF) with the recently proposed paradigm of direct preference optimization (DPO). We focus our attention on the class of loglinear policy parametrization and linear reward functions. In order to compare the two paradigms, we first derive minimax statistical bounds on the suboptimality gap induced by both RLHF and DPO, assuming access to an oracle that exactly solves the optimization problems. We provide a detailed discussion on the relative comparison between the two paradigms, simultaneously taking into account the sample size, policy and reward class dimensions, and the regularization temperature. Moreover, we extend our analysis to the approximate optimization setting and derive exponentially decaying convergence rates for both RLHF and DPO. Next, we analyze the setting where the ground-truth reward is not realizable and find that, while RLHF incurs a constant additional error, DPO retains its asymptotically decaying gap by just tuning the temperature accordingly. Finally, we extend our comparison to the Markov decision process setting, where we generalize our results with exact optimization. To the best of our knowledge, we are the first to provide such a comparative analysis for RLHF and DPO.

---

<sup>1</sup>Max Planck Institute for Software Systems, Saarbrücken, Germany <sup>2</sup>University of Warwick, Coventry, UK <sup>3</sup>Independent Researcher, London, UK. Correspondence to: Andi Nika <andinika@mpi-sws.org>.

## 1. Introduction

Learning from human preferences has grown more prominent as we move closer to artificial general intelligence. One of the most effective ways to learn from preferences is through reinforcement learning from human feedback (RLHF), which involves a two-step process of reward learning and regularized policy optimization. The attractiveness of this paradigm lies in its ability to model the reward function based solely on preference data. This makes it highly applicable in numerous practical situations where rewards are not given *a priori* or are challenging to define accurately. Once the reward is modeled, RLHF solves a regularized value function maximization problem to obtain a fine-tuned policy. This paradigm has enjoyed a lot of applications varying from game-playing (Christiano et al., 2017; Warnell et al., 2018; Knox & Stone, 2008; MacGlashan et al., 2017), robotics (Shin et al., 2023; Brown et al., 2019), and training large language models (LLMs) (Ziegler et al., 2019; Nakano et al., 2021; Wu et al., 2021; Ouyang et al., 2022; Stiennon et al., 2020; Glaese et al., 2022; Ramamurthy et al., 2023; Menick et al., 2022; Ganguli et al., 2022; Bai et al., 2022; Gao et al., 2023).

As an alternative to RLHF, Rafailov et al. (2023) have recently proposed direct preference optimization (DPO), an RL-free paradigm to learning from preferences. DPO circumvents the reward modeling phase and directly optimizes the policy parameters based on the preference data. In certain LLM instances, DPO seems to be empirically superior to RLHF, due to its simple optimization framework.

That said, a statistical analysis of the differences between these paradigms is lacking. The sample complexity of RLHF in various settings has already been studied (Zhu et al., 2023; Zhan et al., 2023; Xiong et al., 2023), and there have been some initial attempts at theoretically understanding DPO and its variants (Azar et al., 2023). However, it is unclear when one of the paradigms is better and when these two paradigms are statistically comparable. Motivated by this observation, we initiate a thorough discussion of the theoretical comparison between RLHF and DPO. Specifically, the purpose of this paper is to address the following research questions:

## Reward Model Learning vs. Direct Policy Optimization

		RLHF	DPO
Realizable Rewards	Exact Optimization	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \tilde{\Theta}\left(\sqrt{\frac{d_R}{n}}\right)$	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \Theta\left(\frac{d_P}{\beta n}\right)$
	Approximate Optimization	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \tilde{\Theta}\left(\sqrt{\frac{d_R}{n}}\right) + O\left((1 - \frac{1}{n})^t + \frac{e^{-t}}{\beta}\right)$	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \Theta\left(\frac{d_P}{\beta n}\right) + O\left(\frac{1}{\beta}\left(1 - \frac{\beta}{n}\right)^t\right)$
Non-realizable Rewards	Exact Optimization	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \tilde{\Theta}\left(\sqrt{\frac{d_R}{n}}\right) + O(\epsilon_{\text{app}})$	$O(\beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \mu)) + \Theta\left(\frac{d_P}{\beta n}\right) + O(\beta D_{\text{KL}}(\pi_{\theta^*} \pi_{r^*}^{\text{opt}}))$

Table 1: A presentation of the bounds on the suboptimality gap for RLHF and DPO. The first two rows present bounds under the realizable reward assumption in the exact and approximate optimization frameworks; the last row presents the bounds when the ground-truth reward function is not realizable. Here,  $\pi_{r^*}^{\text{opt}}$  denotes an optimal policy with respect to the ground-truth reward function  $r^*$ ,  $\pi_{r^*}$  denotes an optimal regularized policy, and  $\pi_{\theta^*}$  denotes an optimal loglinear regularized policy. Moreover,  $\beta$  denotes the regularization temperature,  $D_{\text{KL}}$  denotes the KL divergence,  $d_R$  denotes the reward dimension, and  $d_P$  denotes the policy dimension. Finally,  $n$  denotes the sample size,  $t$  denotes the optimization steps for the approximate setting,  $\epsilon_{\text{app}}$  denotes the reward mismatch coefficient and  $\tilde{\Theta}$  hides any log factors.

*What are the statistical guarantees of RLHF relative to those of DPO? What conditions benefit one as opposed to the other?*

As DPO does not learn a reward model, but directly optimizes over the policy space, a dependence on the policy dimensionality  $d_P$  is expected. On the other hand, RLHF’s performance evidently implies some dependence on reward dimensionality  $d_R$  due to its reward learning phase. Does this imply a discrepancy in the statistical bounds of these paradigms when the reward and policy dimensions are different? Moreover, what can be said about the dependency of the bounds on the sample size  $n$  or the regularization temperature  $\beta$ ?

We address these questions in the following setting: finite spaces, Bradley-Terry preference model, linear rewards and loglinear policies. We first study the exact optimization setting and derive bounds on both RLHF and DPO. Then, we proceed to derive fast convergence rates for a modified version of the policy gradient for RLHF, and gradient descent for DPO. Next, we discuss some implications of our bounds when the reward function is not fully realizable. We close our paper by extending our comparative analysis to deterministic Markov decision processes. Our contributions are summarized below; see Table 1 for explicit bounds.

- First, we derive minimax bounds on the suboptimality gap induced by RLHF and DPO in the exact optimization setting by leveraging smoothness and strong convexity properties. We show that when the optimal regularized policy is loglinear and the reward function is linear, RLHF is  $\tilde{\Theta}(\sqrt{d_R/n})$ -close to its objective, while DPO is  $\tilde{\Theta}(d_P/(\beta n))$ -close. These results emphasize the comparison of the two paradigms in terms of the reward and policy dimensions when setting  $\beta = \Theta(\sqrt{d_P/n})$  for DPO.

- Furthermore, we study the convergence rates of a version of the natural policy gradient for RLHF and gradient descent for DPO. Motivated by recent fast convergence results for entropy regularized RL with tabular softmax policies, we derive  $O(e^{-t}/\beta)$  convergence rates in  $t$  iterations for a version of the natural policy gradient for RLHF. Moreover, for gradient descent we are able to show  $O((1/\beta)(1 - \beta/n)^t)$  convergence rates by using the fact that the DPO loss function satisfies the *PL condition* (Karimi et al., 2016). These results replicate the implications from the exact optimization setting on the difference in terms of reward and policy dimensions.
- We also consider the case where the ground-truth reward function is not realizable and its best linear fit is  $\epsilon_{\text{app}}$ -close to it. We show that, while RLHF incurs an additional constant term on the suboptimality gap, DPO’s dependence on the additional term can be controlled by setting the regularization temperature accordingly.
- Finally, we extend our comparison to deterministic Markov decision processes by proposing a new formulation of the DPO objective for this setting and then generalizing our results. The main motivation for this extension is that, arguably, the difference in reward and policy dimensions in this setting is higher.

## 2. Related Work

**Learning from pairwise comparisons.** In the context of RL, the problem of learning from pairwise comparisons has been studied thoroughly in the bandit setting, where the problem is known as the *dueling bandit* problem (Yue et al., 2009; Faury et al., 2020; Ailon et al., 2014; Gajane et al., 2015; Komiyama et al., 2015; Zoghi et al., 2014; Saha & Gopalan, 2019; Saha & Krishnamurthy, 2022). For the case

of dueling RL for linear MDPs, Saha et al. (2023) propose an algorithm that satisfies tight regret guarantees, while Chen et al. (2022) extend this formulation to the MDPs with general function approximation. Finally, Chatterji et al. (2021) consider a more general setting in which the trajectory-based feedback is generated from a generalized linear model, and they propose variants of optimistic algorithms for online RL.

In this paper, we consider the offline setting, where Zhu et al. (2023) and Zhan et al. (2023) have already provided statistical bounds on pessimistic RLHF for direct value maximization. Our focus, however, is on the regularized value maximization problem. While pessimism mitigates poor coverage for the setting considered by Zhu et al. (2023) and Zhan et al. (2023), for regularized RLHF and, as a consequence, for DPO, this issue remains and is captured by the coverage coefficients which we define with respect to both reward and policy features. The aim of this paper, however, is not to mitigate these issues, but to thoroughly analyse the statistical guarantees of the existing prominent paradigms of learning from human preferences and derive insightful results that shed light on their relative performances.

**Direct preference optimization.** In recent works, the RL-free fine-tuning paradigm of direct preference optimization (DPO) has gained popularity (Rafailov et al., 2023; An et al., 2023; Azar et al., 2023; Wang et al., 2023). Its original formulation was proposed for the contextual bandit setting. Hejna et al. (2023) propose an extension of DPO to MDPs under the assumption that the preferences depend on the advantage function of the optimal policy. While we also provide an extension of the DPO formulation for MDPs in Section 7, our primary focus is on a comparative analysis between the two paradigms in the contextual bandit setting.

**Offline RL.** In recent years, there has been a significant surge in interest towards offline RL, with an extensive literature both in the empirical front (Jaques et al., 2019; Laroché et al., 2019; Fujimoto et al., 2019; Kumar et al., 2020; Agarwal et al., 2020; Kidambi et al., 2020) and the theoretical one (Jin et al., 2021; Xie et al., 2021; Rashidinejad et al., 2021; Uehara & Sun, 2021; Zanette et al., 2021). While the focus of this line of work is on the traditional reward-based offline RL, our problem is derived from a combination of reward-learning from pairwise feedback and KL-regularized offline RL based on it.

### 3. Formal Setting

This section presents the background material that will be used throughout the paper. We will use a notation similar to (Rafailov et al., 2023) and (Azar et al., 2023).

**Notation.** Let  $\langle u, v \rangle = u^\top v$  denote the inner product between vectors  $u$  and  $v$ . The trace of a matrix  $A$  is denoted by  $\text{tr}(A)$  and its pseudo-inverse by  $A^\dagger$ . Moreover,  $\Delta(\mathcal{X})$  denotes the set of distributions over the finite set  $\mathcal{X}$  and  $\|v\|_M = \sqrt{v^\top M v}$  denotes the seminorm of vector  $v$  with respect to  $M$ . Finally,  $\text{proj}_A(v)$  denotes the projection of vector  $v$  onto set  $A$  and  $\tilde{\Theta}(\cdot)$  hides any log factors.

#### 3.1. Preliminaries

Let  $\mathcal{X}$  be a finite set of contexts with cardinality  $X$  and  $\mathcal{Y}$  be a finite set of actions with cardinality  $Y$ . Fix  $\rho \in \Delta(\mathcal{X})$  as an initial distribution over contexts and let  $r : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$  be a reward function. We consider a class of linear reward functions defined below.

**Definition 3.1 (Linear reward function class).** Let  $\phi$  be a  $d_R$ -dimensional feature mapping with  $\max_{x,y} \|\phi(x,y)\|_2 \leq 1$  and let  $F > 0$ . We consider the following class of linear reward functions:

$$\mathcal{F} = \left\{ r_\omega \in [0, 1]^{XY} : r_\omega(x, y) = \omega^\top \phi(x, y), \right. \\ \left. \forall (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ where } \omega \in \mathbb{R}^{d_R} \text{ and } \|\omega\|_2 \leq F \right\}.$$

Given  $x \in \mathcal{X}$ , a policy  $\pi(\cdot|x) \in \Delta(\mathcal{Y})$  is a distribution over actions. Throughout the paper, we will consider the loglinear class of policies, defined as follows.

**Definition 3.2 (Loglinear policy class).** Let  $\psi$  be a  $d_P$ -dimensional feature mapping with  $\max_{x,y} \|\psi(x,y)\|_2 \leq 1$  and let  $B > 0$ . We consider the following class of loglinear policies:

$$\Pi = \left\{ \pi_\theta : \pi_\theta(y|x) = \frac{\exp(\theta^\top \psi(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y'))}, \right. \\ \left. \forall (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ where } \theta \in \mathbb{R}^{d_P} \text{ and } \|\theta\|_2 \leq B \right\}.$$

Given policy  $\pi$ , the value function of  $\pi$  with respect to reward function  $r$  and context distribution  $\rho$  is defined as

$$V_r^\pi(\rho) = \sum_x \rho(x) \sum_y \pi(y|x) r(x, y).$$

#### 3.2. Offline Learning from Human Preferences

Let  $\mu$  be a reference policy fixed throughout the paper, and  $\beta > 0$  be a regularization parameter. Let us define the KL-regularized objective as

$$\mathcal{V}_r^\pi(\rho) = V_r^\pi(\rho) - \beta D_{\text{KL}}(\pi || \mu),$$

where  $D_{\text{KL}}(\pi || \mu) = \sum_x \rho(x) \sum_y \pi(y|x) \log \frac{\pi(y|x)}{\mu(y|x)}$ .

We assume access to the dataset  $\mathcal{D}_n = \{(x_i, y_i^w, y_i^l)\}_{i=1}^n$ , where  $y_i^w$  denotes the preferred action over  $y_i^l$ . In this paper,

we will assume that the distribution of human preferences follows the Bradley-Terry (BT) model (Bradley & Terry, 1952), which we formally state below.

**Definition 3.3** (*Bradley-Terry preference model*). There exists a latent reward function  $r^*$  and a probability law  $P^*$  such that, for every tuple  $(x, y^w, y^l)$ , we have

$$P^*(y^w \succ y^l | x) = \frac{\exp(r^*(x, y^w))}{\exp(r^*(x, y^w)) + \exp(r^*(x, y^l))},$$

where  $y^w \succ y^l$  denotes  $y^w$  being preferred over  $y^l$ .

The latent reward function  $r^*$  will be fixed throughout the paper as the ground-truth reward function.

### 3.3. Reinforcement Learning from Human Feedback

We consider the reinforcement learning from human feedback (RLHF) paradigm as formulated in (Ziegler et al., 2019). Having access to preference dataset  $\mathcal{D}_n$  and a fixed reference policy  $\mu$ , RLHF proceeds in two phases: the reward learning phase and the final KL-regularized reinforcement learning phase.

For the reward learning phase, RLHF estimates  $r^*$  by applying maximum likelihood estimation (MLE) to the dataset  $\mathcal{D}_n$ . The MLE optimization problem can be written as

$$\min_r \mathcal{L}_{\text{RLHF}}^r(\mathcal{D}_n) := -\mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_n} \left[ \log \left( \sigma(r(x^w, y^w) - r(x^l, y^l)) \right) \right], \quad (\text{P1.1})$$

where  $\sigma(z) = 1/(1 + \exp(-z))$  denotes the sigmoid function. Let  $\hat{r}$  denote the solution of Problem (P1.1).

The final phase of RLHF consists of maximizing the KL-regularized objective with respect to  $\hat{r}$  by solving

$$\max_{\pi} \mathcal{V}_{\hat{r}}^{\pi}(\mathcal{D}_n) := \mathbb{E}_{\substack{x \sim \mathcal{D}_n \\ y \sim \pi(\cdot|x)}} \left[ \hat{r}(x, y) - \beta \log \frac{\pi(y|x)}{\mu(y|x)} \right]. \quad (\text{P1.2})$$

### 3.4. Direct Preference Optimization

Recently, alternative paradigms to RLHF have been studied. In particular, Rafailov et al. (2023) introduced direct preference optimization (DPO), a new fine-tuning paradigm that directly optimizes the policy parameters instead of going through the reward modeling phase. Their key observation is that the latent reward can be expressed in terms of its optimal policy and the reference policy. This yields a loss function that is directly defined in terms of the preference data.

Formally, Rafailov et al. (2023) show that there exists a policy  $\pi$  that maximizes the KL-regularized objective, for which we have

$$r^*(x, y) = \beta \log \frac{\pi(y|x)}{\mu(y|x)} + \beta \log Z(x) \quad (1)$$

for every  $(x, y)$ , where  $Z(x)$  denotes the partition function  $\sum_y \mu(y|x) \exp(r^*(x, y)/\beta)$ . A new objective is then derived, which directly depends on the policy. Given preference dataset  $\mathcal{D}_n$ , this objective leads to the following optimization problem:

$$\min_{\pi} \mathcal{L}_{\text{DPO}}^{\pi}(\mathcal{D}_n) := -\mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_n} \left[ \log \left( \sigma \left( \beta \log \frac{\pi(y^w|x)}{\mu(y^w|x)} - \beta \log \frac{\pi(y^l|x)}{\mu(y^l|x)} \right) \right) \right]. \quad (\text{P2})$$

As it turns out, this elegant approach yields practical benefits. However, it is unclear whether these benefits can be theoretically justified. In this paper, we will provide a comparative analysis of both RLHF and DPO in different settings. Next, we define a unified metric of performance to compare these paradigms.

### 3.5. Performance Metric

Given a reward function  $r$ , let  $\pi_r^{\text{opt}} \in \arg \max_{\pi} V_r^{\pi}(\rho)$  denote an optimal policy for  $V_r^{\pi}(\rho)$  and  $V_r^{\text{opt}}(\rho)$  the optimal value. For a given policy  $\pi$ , we define the suboptimality gap of  $\pi$  as

$$G(\pi) = V_{r^*}^{\text{opt}}(\rho) - V_r^{\pi}(\rho).$$

$G(\pi)$  captures how well a policy is performing w.r.t. the ground-truth reward function  $r^*$ . In this paper, we will use the suboptimality gap  $G$  as our unified measure of performance when comparing the two paradigms.

We note that RLHF and DPO are designed to minimize regularized objectives  $\mathcal{V}$  instead of optimizing the value function  $V$  (see Sections 3.3 and 3.4). In order to rigorously analyze the differences between RLHF and DPO, we need to establish some additional notation. Let  $\pi_r^* \in \arg \max_{\pi} \mathcal{V}_r^{\pi}(\rho)$  denote a regularized optimal policy with respect to  $r$  and  $\mathcal{V}_r^*(\rho)$  denote the optimal regularized value. Analogous to  $G$ , we define the regularized suboptimality gap of  $\pi$  as

$$\mathcal{G}(\pi) = \mathcal{V}_{r^*}^*(\rho) - \mathcal{V}_r^{\pi}(\rho) = (V_{r^*}^{\pi_r^*}(\rho) - V_r^{\pi}(\rho)) - \beta(D_{KL}(\pi_{r^*}^* || \mu) - D_{KL}(\pi || \mu)).$$

As RLHF and DPO are designed to minimize  $\mathcal{G}(\pi)$ , both will incur an additional term on their bounds when comparing their bounds w.r.t.  $G(\pi)$ . For this reason, we will formally define this discrepancy term as  $D(\pi) = G(\pi) - \mathcal{G}(\pi)$  and discuss it further in Section 4. Next, we proceed to provide a comparative analysis for these paradigms, starting with the exact optimization setting when the ground-truth reward is realizable. We note that any log terms are omitted for clarity of presentation. All proofs and related discussions can be found in the Appendix.

## 4. Realizable Rewards: Exact Optimization

In this section, we analyze the statistical differences in performance between RLHF and DPO in the exact optimization setting. We assume throughout the section that the ground-truth reward function is linear and realizable, i.e.  $r^* \in \mathcal{F}$ . Moreover, we assume a loglinear regularized optimal policy exists, i.e.  $\pi_{r^*}^* \in \Pi$ . Note that, for linear reward function,  $\mathcal{L}_{\text{RLHF}}^r(\mathcal{D}_n)$  can be equivalently written as

$$-\mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_n} \left[ \log \sigma \left( \omega^\top (\phi(x, y^w) - \phi(x, y^l)) \right) \right]. \quad (2)$$

Moreover, for loglinear policies,  $\mathcal{L}_{\text{DPO}}^\pi(\mathcal{D}_n)$  can be equivalently written as

$$-\mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_n} \left[ \log \left( \sigma \left( \beta \theta^\top (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l) \right) \right) \right], \quad (3)$$

where  $J(x, y^w, y^l) = \log(\mu(y^w|x)/\mu(y^l|x))$ .

We will denote the losses for the RLHF reward learning phase and DPO as  $\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$  and  $\mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)$ , respectively. Let  $r_{\hat{\omega}}$  denote the reward estimate and let  $\pi_{\hat{\theta}}$  denote the policy learned by RLHF. Moreover, let  $\pi_{\tilde{\theta}}$  denote the policy learned by DPO. Formally, we assume that RLHF has access to an oracle that exactly solves both optimization problems and returns  $r_{\hat{\omega}} \in \arg \min \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$  and  $\pi_{\hat{\theta}} \in \arg \max_{\theta} \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\mathcal{D}_n)$ . Similarly, for DPO we assume that the oracle returns  $\pi_{\tilde{\theta}} \in \arg \min_{\theta} \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)$ .

### 4.1. Theoretical Results

Before stating our main results of this section, we will need to define the covering numbers with respect to  $\mathcal{D}_n$ . Let the sample covariance matrix with respect to the reward features be defined as

$$\Sigma_{\mathcal{D}_n, R} = \frac{1}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \bar{\phi}(x, y^w, y^l) \bar{\phi}(x, y^w, y^l)^\top,$$

where  $\bar{\phi}(x, y^w, y^l) = (\phi(x, y^w) - \phi(x, y^l))$ , for every  $(x, y^w, y^l) \in \mathcal{D}_n$ . Fix  $\lambda > 0$  and let  $\Lambda_R = \left\| (\Sigma_{\mathcal{D}_n, R} + \lambda I)^{-1/2} \right\|_2$  be the reward covering number.

Similarly, let the sample covariance matrix of the policy features be defined as

$$\Sigma_{\mathcal{D}_n, P} = \frac{1}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \bar{\psi}(x, y^w, y^l) \bar{\psi}(x, y^w, y^l)^\top,$$

where  $\bar{\psi}(x, y^w, y^l) = (\psi(x, y^w) - \psi(x, y^l))$ , for every  $(x, y^w, y^l) \in \mathcal{D}_n$ . The policy covering number of  $\mathcal{D}_n$  is  $\Lambda_P = \left\| (\Sigma_{\mathcal{D}_n, P} + \lambda I)^{-1/2} \right\|_2$ . Our bounds will depend on

these quantities. We will start with minimax bounds on the suboptimality for RLHF in the above-mentioned setting.

**Theorem 4.1.** *Let  $\delta > 0$ . Assume that  $r^* \in \mathcal{F}$ . Then, with probability at least  $1 - \delta$ , the suboptimality gap incurred by RLHF is*

$$G(\pi_{\hat{\theta}}) = \Theta \left( \Lambda_R \sqrt{\frac{d_R + \log(6/\delta)}{S_R^2 n}} + \lambda F^2 \right) + D(\pi_{\hat{\theta}}),$$

where  $S_R = 1 / (2 + \exp(-2F) + \exp(2F))$ .

*Proof sketch.* The proof follows from splitting the gap into sub-gaps which we can bound directly, and results from (Zhu et al., 2023).  $\square$

Next, we will consider the suboptimality gap induced by DPO. First, note that Equation (3) for loglinear policies essentially becomes a logistic regression problem in  $d_P + 1$  dimensions, by adding a dummy variable to  $\theta$  that corresponds to  $J(x, y^w, y^l)$ . Our goal is to use loglinearity to derive smoothness properties for logistic regression so that we can obtain minimax bounds – without this assumption on the policy class,  $\mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)$  may not satisfy these properties.<sup>1</sup>

Moreover, note that Equation (1) relates  $r^*$  to one of the regularized optimal policies  $\pi_{r^*}^*$ . Obviously, this does imply that  $\pi_{r^*}^*$  is loglinear. Nevertheless, the following lemma states that that is the case for linear rewards. Let  $\Phi \in \mathbb{R}^{d_R \times XY} = [\phi(x, y)]_{(x, y) \in \mathcal{X} \times \mathcal{Y}}$  and  $\Psi \in \mathbb{R}^{d_P \times XY} = [\psi(x, y)]_{(x, y) \in \mathcal{X} \times \mathcal{Y}}$  denote the reward feature and policy feature matrices, respectively.

**Lemma 4.1.** *Assume that  $r^* \in \mathcal{F}$  and  $\mu \in \Pi$ . Furthermore, assume that the column space of  $\Phi$  is a subspace of the column space of  $\Psi$ . Then, there exists  $\theta^* \in \Theta$ , such that  $\pi_{\theta^*} \in \arg \max_{\pi} \mathcal{V}_{r^*}^\pi(\rho)$  and  $r^*(x, y) = \beta \log(\pi_{\theta^*}(y|x)/\mu(y|x)) + \beta \log Z(x)$ .*

With these observations in place, we are now ready to state the minimax bounds on the suboptimality for DPO.

**Theorem 4.2.** *Let  $\delta > 0$  and  $\beta > 0$ . Assume that  $r^* \in \mathcal{F}$ ,  $\mu \in \Pi$ , and that the condition of Lemma 4.1 is satisfied. Let  $n \geq O\left(\text{tr}(\Sigma_{\mathcal{D}_n, P}^\dagger)/(\beta B^2)\right)$ . Then, with probability at least  $1 - \delta$ , the suboptimality gap of DPO is*

$$G(\pi_{\tilde{\theta}}) = D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda_P(d_P + 1)}{\beta n} + \beta \lambda \Lambda_P B^2 \right).$$

*Proof sketch.* We start by splitting the suboptimality gap and focus on the  $\mathcal{G}(\pi_{\tilde{\theta}})$  term. Here, we utilize the expression for the ground-truth reward as  $r^*(x, y) = \beta \log \pi_{\theta^*}(y|x) - \beta \log \mu(y|x) + Z_{\tilde{\theta}}(x)$ , where  $Z_{\tilde{\theta}}(x)$  denotes the partition function with respect to  $\pi_{\tilde{\theta}}$ . Then, using the fact that the policies are loglinear we further expand and reduce the whole

<sup>1</sup>Lemma J.6 shows this loss is not smooth for tabular settings.

gap in terms of the log differences. Next, we utilize the smoothness and Lipschitzness of the log-sum-exp function to finally obtain the upper bounds. For the lower bound, we construct an example where the policy feature matrix  $\Psi$  is full rank, and show that the log-sum-exp function becomes strongly convex. This finally leads to the stated bounds.  $\square$

Before discussing the implications of our results, let us say a few words on the regularization gap  $D(\pi)$  for RLHF and DPO. Let  $\pi_{\theta^*}$  denote an optimal loglinear regularized policy. A characterization of  $D(\pi)$  is given as follows.

**Lemma 4.2.** *For any  $\theta$ , we have that  $\beta D_{\text{KL}}(\pi_{\theta^*} \parallel \mu) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \mu) \leq D(\pi_{\theta})$  and  $D(\pi_{\theta}) \leq \beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \parallel \mu) - \beta D_{\text{KL}}(\pi_{\theta} \parallel \mu)$ .*

Furthermore, for DPO, this quantity can be upper bounded in terms of  $\pi_{r^*}^{\text{opt}}$  and  $\pi_{\theta^*}$  as follows.

**Lemma 4.3.** *Given  $\delta > 0$ , with probability at least  $1 - \delta$ , we have  $D(\pi_{\tilde{\theta}}) \leq \beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \parallel \mu) - \beta D_{\text{KL}}(\pi_{\theta^*} \parallel \mu) + \tilde{O}(d_P/n^{3/2})$ .*

In general, it is known that the KL divergence may not be upper-bounded. However, assuming that optimal policy  $\pi_{r^*}^{\text{opt}}$ , optimal regularized policy  $\pi_{\theta^*}$ , and sampling policy  $\mu$  are not far away from each other, then these  $D_{\text{KL}}(\cdot)$  quantities would not be too large.

## 4.2. Comparative Analysis

In this section, we will provide some insights into the implications of our theoretical results. For the purpose of this section, we will focus our attention on the problem-dependent parameters and ignore the quantity  $D(\pi)$ .

**The role of dimensionality.** Note that RLHF has  $\tilde{\Theta}(\sqrt{d_R})$  dependence on the reward dimension, while DPO has  $\Theta(d_P)$  dependence on the policy dimension. When  $d_R = d_P$  and the sample size is small, RLHF seems to statistically outperform DPO. Any setting where  $d_R \ll d_P$  makes this difference more apparent. In Section 7, we will discuss an extension of our analysis to a setting where the reward dimension can be much smaller than the policy dimension in practice.

**The role of sample size.** Next, we take into consideration the sample size. Note that DPO’s bounds depend on  $n$  being large enough (cf. Theorem 4.2). Assume everything else constant and  $d = d_R = d_P$ . If the  $D(\pi)$  terms are similar for both paradigms, then, for large sample sizes such that  $n \gg d$ , DPO seems to outperform RLHF asymptotically. Whenever  $n < d$  (which is usually the case for large language models), RLHF has a smaller suboptimality gap.

**The role of  $\beta$ .** Finally, we discuss the role of the temperature  $\beta$  on the bounds. First, note that RLHF can effectively set  $\beta = 0$  to annihilate the effect that  $D(\pi)$  has on its bounds. On the other hand, DPO cannot set  $\beta$  to 0 due to a

disproportional dependence of its bounds on it. Thus, the optimal choice of  $\beta$  for DPO is  $\beta = \Theta(\sqrt{d_P/n})$ , yielding  $\Theta(\sqrt{d_P/n})$  bounds and matching the order of  $n$  in the bounds of RLHF. For such a value of  $\beta$ , the same implications hold – the main difference between both settings is in terms of the differences between the reward and policy parameter dimensions.

## 5. Realizable Rewards: Approximate Optimization

In this section, we shift our focus to the approximate setting, where access to oracles is not given. Here, both paradigms have to approximately solve their estimation problems based on the given data. Similar to the previous section, we assume throughout this section that the ground-truth reward function  $r^*$  is linear and realizable in  $\mathcal{F}$ , and that there exists a loglinear regularized policy  $\pi_{r^*} \in \Pi$ . Moreover, we assume that, for every data tuple  $(x, y^w, y^l) \in \mathcal{D}_n$ , we have  $\phi(x, y^w) \neq \phi(x, y^l)$  and  $\psi(x, y^w) \neq \psi(x, y^l)$ .

### 5.1. Theoretical Results

Let us start with the reward learning phase. Recall the definition of the loss  $\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$  for MLE, as defined in Section 4. Let  $\omega_0$  be initialized randomly, and let

$$\omega_{t+1} = \underset{\omega: \|\omega\|_2 \leq F}{\text{proj}} (\omega_t - \eta \nabla_\omega \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)), \quad (4)$$

for any iterate  $t \geq 0$ , where  $\eta$  denotes the learning rate. Let  $\omega_{\mathcal{D}_n}^* \in \arg \max \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$ . The first result of this section provides fast convergence rates of gradient descent for the reward learning phase of RLHF.

**Theorem 5.1.** *For every  $t \geq 0$ , the gradient descent procedure (4) with learning rate  $\eta = 1/\exp(2F)$  satisfies*

$$\|\omega_t - \omega_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, R}}^2 \leq O\left(1 - \frac{1}{n}\right)^t.$$

*Proof sketch.* We begin by showing Lipschitzness and smoothness of  $\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$  with respect to  $\omega$ . Then, we show that the PL condition (Karimi et al., 2016) is enough to guarantee fast convergence of projected gradient descent by showing that such a condition implies that  $\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$  also satisfies the proximal PL condition (Karimi et al., 2016) when its domain is restricted to a ball. The result follows by applying the convexity of  $\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)$ .  $\square$

Next, we discuss the policy optimization phase of RLHF. Let  $r_\infty$  denote the reward estimated from the previous phase. Initialize  $\theta_0 \in \mathbb{R}^{d_P}$  with  $\|\theta_0\|_2 \leq B$ . For any  $t \geq 0$ , let

$$\theta_{t+1} = \theta_t + \eta' (\Psi_n \Psi_n^\top)^\dagger \nabla_\theta \mathcal{V}_{r_\infty}^{\pi_\theta}(\mathcal{D}_n), \quad (5)$$

where  $\eta' > 0$  is the learning rate and  $\Psi_n = [\psi(x, y)]_{x \in \mathcal{D}_n, y \in \mathcal{Y}}$  denotes the sample feature matrix.<sup>2</sup> Then, the following result holds.

**Theorem 5.2.** *Let  $\delta > 0$ . Assume that  $\Psi_n$  has full column rank. Then, with probability at least  $1 - \delta$ , for every  $t \geq 1$ , update rule (5) with learning rate  $\eta' \leq n/\beta$  satisfies*

$$\mathcal{V}_{r_{\tilde{\omega}}}^*(\mathcal{D}_n) - \mathcal{V}_{r_{\tilde{\omega}}}^{\pi_{\tilde{\theta}^t}}(\mathcal{D}_n) \leq O\left(\frac{1}{\beta} \exp(-(t-1))\right).$$

*Proof sketch.* After deriving a naive gradient update rule in matrix notation, we examine the conditions needed to obtain fast convergence rates for our setting. As a consequence, we design our gradient update which resembles natural policy gradient for loglinear policies – the matrix  $(\Psi_n \Psi_n^\top)^\dagger$  captures the gradient information for this case. Then, we use similar techniques to those in (Mei et al., 2020) and use the fact that  $\Psi_n$  is full rank, to finally obtain the desired bounds.  $\square$

*Remark 5.1.* It is important to emphasize that the assumption on  $\Psi_n$  is not restrictive. Indeed, given a preference dataset, it is always possible to construct alternative feature representations that satisfy the assumption with respect to the given data (e.g. different LLMs use different encoding methods while being fine-tuned using the same data).

*Remark 5.2.* Before going to our next result, it is important to clarify the double usage of  $\mathcal{D}_n$  for both the reward learning and policy optimization phases. Our theoretical guarantees are based on the data being independently generated in these phases. Thus, the standard approach is to split the data into two batches for both purposes. Note that both batches would still be  $O(n)$  in size and the dependence of the results on  $n$  would not change. We use the same  $\mathcal{D}_n$  for both phases for simplicity of presentation.

Next, we provide convergence results of gradient descent for DPO with loglinear policies. Let  $\theta_0$  be initialized randomly, and let

$$\theta_{t+1} = \text{proj}_{\theta: \|\theta\| \leq B} (\theta_t - \eta'' \nabla_{\theta} \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)), \quad (6)$$

for any iterate  $t \geq 0$ , where  $\eta''$  denotes the learning rate. Let  $\theta_{\mathcal{D}_n}^* \in \arg \max \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)$ . Then, we have the following result.

**Theorem 5.3.** *For every  $t \geq 0$ , the gradient descent procedure (6) with learning rate  $\eta'' = O(1/\beta^2)$  satisfies*

$$\|\theta_t - \theta_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, P}}^2 \leq O\left(\frac{1}{\beta} \left(1 - \frac{\beta}{n}\right)^t\right).$$

We have omitted the dependence on the absolute constants that are irrelevant to our discussion – see Appendix F for a detailed expression of the hidden constants.

<sup>2</sup>We only need  $\|\theta_0\|_2 \leq B$  for RLHF, since its bounds do not depend on  $\|\theta_t\|$ , for  $t \geq 1$ . Thus, we do not need projection.

## 5.2. Comparative Analysis

The regularized suboptimality gap for RLHF is

$$\mathcal{G}(\pi_{\tilde{\theta}}) \leq \Theta\left(\sqrt{\frac{d_R}{n}}\right) + O\left(\left(1 - \frac{1}{n}\right)^t + \frac{\exp(-t)}{\beta}\right)$$

and the regularized suboptimality gap for DPO is

$$\mathcal{G}(\pi_{\tilde{\theta}}) \leq \Theta\left(\frac{d_P}{\beta n}\right) + O\left(\frac{1}{\beta} \left(1 - \frac{\beta}{n}\right)^t\right).$$

Both of these paradigms satisfy exponential convergence rates, thus, the main implications of the discussion of Section 4.2 hold in this setting as well if  $\beta$  is to be set as constant. If  $\beta$  is to be tuned for DPO, it cannot be made arbitrarily small or large as observed in Section 4.2 – DPO’s overall bounds disproportionately depend on the parameter  $\beta$ . Even so, setting  $\beta$  to its optimal value of  $\Theta(1/\sqrt{n})$  for the exact optimization setting would not affect the convergence rate of gradient descent for DPO.

## 6. Non-realizable Rewards: Exact Optimization

In this section, we consider the case when the ground-truth reward function  $r^*$  does not belong to the linear class  $\mathcal{F}$  (see Definition 3.1). We again assume that there exists an optimal  $\pi_{r^*}$  regularized policy that belongs to the loglinear class  $\Pi$  (see Definition 3.2) for some  $\theta^*$ . We will capture the mismatch between the reward function and its best linear approximation in  $\mathcal{F}$  by the following condition.

**Assumption 6.1** (*Non-realizability of the ground-truth reward*). There exists  $r_{\omega^*} \in \mathcal{F}$  with parameter  $\omega^*$  such that  $\|r^* - r_{\omega^*}\|_{\infty} \leq \epsilon_{\text{app}}$ , where  $\epsilon_{\text{app}} > 0$  denotes the mismatch coefficient.

As we will see, the mismatch coefficient will appear linearly in the RLHF bounds on the gap as an additional constant that cannot be improved by increasing the data size.

### 6.1. Theoretical Results

We begin with the RLHF result, which can be derived from Theorem 4.1.

**Theorem 6.1.** *Let  $\delta > 0$ . Suppose that Assumption 6.1 holds. Then, with probability at least  $1 - \delta$ , we have  $G(\pi_{\tilde{\theta}}) \leq D(\pi_{\tilde{\theta}}) + \tilde{\Theta}\left(\Lambda_R \sqrt{d_R/n}\right) + 2\epsilon_{\text{app}}$ .*

We can directly obtain a similar dependence on  $\epsilon_{\text{app}}$  for DPO. In addition to that, we can also obtain alternative bounds that can be controlled by  $\beta$  as follows.

**Theorem 6.2.** *Let  $\delta > 0$ . Suppose that Assumption 6.1 and the condition of Lemma 4.1 hold. Then, with probability at least  $1 - \delta$ , we have  $G(\pi_{\tilde{\theta}}) \leq D(\pi_{\tilde{\theta}}) + \Theta(\Lambda_P d_P / (\beta n)) + \min\{2\epsilon_{\text{app}}, O(\beta D_{\text{KL}}(\pi_{\theta^*} \|\pi^*))\}$ .*

## 6.2. Comparative Analysis

The key observation to be made in this section is the discrepancy of the bounds in terms of the reward mismatch coefficient. RLHF does not use MLE for the policy parameter estimation, but first learns a reward model. Thus, it cannot bypass the error coming from the reward unrealizability. For DPO, note that, if we set  $\beta = O(1/\sqrt{n})$ , its bounds improves asymptotically with  $n$ , assuming that  $D_{\text{KL}}(\pi_{\theta^*} \|\pi_{r^*}^*)$  is bounded (recall that  $\pi_{r^*}^*$  denotes the optimal regularized policy and  $\pi_{\theta^*}$  denotes its best loglinear fit). This setting benefits DPO as it is designed to bypass the reward function and directly optimize over the policy space.

## 7. Realizable Rewards: Exact Optimization – An Extension to Deterministic MDPs

Up to this point, our discussion was concentrated on the contextual bandit setting, which has been used in the DPO literature for the KL-regularized problem (Rafailov et al., 2023). Now, we focus on a generalization of our comparative analysis to Markov decision processes (MDPs), where contexts are related to each other through transition dynamics.

As mentioned previously, the discrepancy between reward and policy dimensions plays a crucial role in the relative performances of RLHF and DPO. While these dimensions could arguably be similar (or have a small gap) for the contextual bandit setting, that is not necessarily the case in general when extending to MDPs, where the reward dimension can be smaller than the policy dimension. For this section, we assume that the ground-truth reward function  $r^*$  is linear and realizable in  $\mathcal{F}$ .

### 7.1. Preliminaries for Deterministic MDPs

For an MDP,  $\mathcal{X}$  is the set of states and  $\mathcal{Y}$  the set of actions. In particular, we consider deterministic MDPs, with a transition function  $T : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X}$  that provides the next state, given the current state-action.

The value function in infinite-horizon MDPs is given as  $V_r^\pi(x) = \mathbb{E}[\sum_{t \geq 0} \gamma^t r(x_t, y_t) | \rho, \pi]$ , where  $\rho$  denotes the initial state distribution. Given policy  $\pi$ , the occupancy measure of  $\pi$  is given by  $d_\rho^\pi(x, y) = (1 - \gamma) \sum_{t \geq 0} \gamma^t \mathbb{P}(x_t = x, y_t = y | \rho, \pi)$ . We will consider the class of loglinear occupancy measures, as defined next.

**Definition 7.1** (Loglinear occupancy measures class). Let  $\psi'(x, y) \in \mathbb{R}^{d_M}$  denote the feature vector of the pair  $(x, y)$  with  $\max_{x, y} \|\psi'(x, y)\|_2 \leq 1$ , and  $B' > 0$ . We consider the following class of loglinear occupancy measures:

$$\Pi' = \left\{ d_\rho^\theta : d_\rho^\theta(x, y) = \frac{\exp(\theta^\top \psi'(x, y))}{\sum_{x', y'} \exp(\theta^\top \psi'(x', y'))}, \right. \\ \left. \forall (x, y) \in \mathcal{X} \times \mathcal{Y} \text{ where } \theta \in \mathbb{R}^{d_M} \text{ and } \|\theta\|_2 \leq B' \right\}.$$

In this section, we use the  $\mathcal{V}_r^{d_\rho}(\rho)$  notation, instead of  $\mathcal{V}_r^\pi(\rho)$ . Similarly, we use  $G(d_\rho)$  to denote the gap in terms of occupancy measure  $d_\rho$ , and  $D(d_\rho)$  for the difference of the gaps. For a complete discussion, see Appendix H.

In this setting, we are given a dataset  $\mathcal{D}_n = \{(x_{0,i}, \tau_i^w, \tau_i^l)\}_{i=1}^n$ , where  $x_{0,i}$  denotes the initial state of the  $i$ th sample and  $\tau_i^w$  denotes the preferred trajectory  $(x_{0,i}, y_{0,i}^w, x_{1,i}^w, \dots)$  over  $\tau_i^l$ . Analogous to Section 3, we define the Bradley-Terry preference model for two trajectories  $\tau^w$  and  $\tau^l$  as  $P^*(\tau^w \succ \tau^l | x_0) = \sigma(R^*(\tau^w) - R^*(\tau^l))$ , where  $R^*(\tau) = \sum_{t \geq 0} \gamma^t r^*(x_t, y_t)$  is the discounted return, and  $\tau^w \succ \tau^l$  denotes  $\tau^w$  being preferred over  $\tau^l$ .

### 7.2. RLHF and DPO for MDPs

Similar to the contextual bandit setting, the objective for the reward learning phase of RLHF in MDPs with linear rewards can be written as follows:

$$\min_{\omega} \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n) := -\mathbb{E}_{(x_0, \tau^w, \tau^l) \sim \mathcal{D}_n} \left( \log \sigma \left( \omega^\top \left( \sum_{t \geq 0} \gamma^t (\phi(x_t^w, y_t^w) - \phi(x_t^l, y_t^l)) \right) \right) \right). \quad (\text{P3.1})$$

Once we have the estimated reward function  $r_{\hat{\omega}}$ , the objective is to solve the KL-regularized problem. Following previous literature on KL-regularized RL (Nachum et al., 2019; Lee et al., 2021), we formulate the objective in this setting as

$$\max_{\pi} V_{r_{\hat{\omega}}}^\pi(\rho) - \beta D_{\text{KL}}(d_\rho^\pi \| d_\rho^\mu). \quad (\text{P3.2})$$

We will assume throughout this section that we are given access to oracles that exactly solve Problem (P3.1) and (P3.2). *Remark 7.2.* Note that the objective in Problem (P3.2) depends on  $\rho$ , while the objective in Problem (P1.2) depends on  $\mathcal{D}_n$ . This is due to considering occupancy measures instead of policies. We keep our current formulation for ease of presentation and leave its extension to a sample objective formulation for future work.

For the purposes of our comparative analysis, we also need an extension of DPO to MDPs, based on the preference model of Section 7.1. The key difficulty of extending DPO to the MDP setting is that the gradient has a non-linear dependence on the policy. To bypass this issue, we leverage the fact that transitions are deterministic to simplify cumulative differences of the optimal Lagrange multipliers for Problem (P3.2), and obtain the following loss function for DPO:

$$\mathcal{L}_{\text{DPO}}^{d_\rho}(\mathcal{D}_n) = -\mathbb{E}_{(x_0, \tau^w, \tau^l) \sim \mathcal{D}_n} \left[ \log \sigma \left( \beta \sum_{t=0}^{\infty} \gamma^t \left( \log \frac{d_\rho(x_t^w, y_t^w)}{d_\rho^\mu(x_t^w, y_t^w)} - \log \frac{d_\rho(x_t^l, y_t^l)}{d_\rho^\mu(x_t^l, y_t^l)} \right) \right) \right].$$



All derivations are in Appendix H. Next, we generalize the bounds from Section 4 for the above formulations.

### 7.3. Theoretical Results

Analogous to the previous sections, let us define

$$\Sigma'_{\mathcal{D}_n, R} = \frac{1}{n} \sum_{(x_0, \tau^w, \tau^l) \in \mathcal{D}_n} \bar{\phi}'(x_0, \tau^w, \tau^l) \bar{\phi}'(x_0, \tau^w, \tau^l)^\top,$$

where  $\bar{\phi}'(x_0, \tau^w, \tau^l) = \sum_{t \geq 0} \gamma^t (\phi(x_t, y_t^w) - \phi(x_t, y_t^l))$ . For  $\lambda > 0$ , define  $\Lambda'_R = \|(\Sigma'_{\mathcal{D}_n, R} + \lambda I)^{-1/2}\|_2$ . Similarly, let the sample covariance matrix with respect to occupancy measure features be defined as

$$\Sigma'_{\mathcal{D}_n, M} = \frac{1}{n} \sum_{(x_0, \tau^w, \tau^l) \in \mathcal{D}_n} \bar{\psi}'(x_0, \tau^w, \tau^l) \bar{\psi}'(x_0, \tau^w, \tau^l)^\top,$$

where  $\bar{\psi}'(x_0, \tau^w, \tau^l) = \sum_{t \geq 0} \gamma^t ((\psi'(x_t, y_t^w) - \psi'(x_t, y_t^l)))$ . Let  $\Lambda'_M = \|(\Sigma'_{\mathcal{D}_n, M} + \lambda I)^{-1/2}\|_2$ . For RLHF with exact optimization, it is straightforward to extend our previous bounds as follows.

**Theorem 7.1.** *Let  $\delta > 0$ . Assume that the policy learning phase yields  $\hat{\theta} \in \arg \max_{\theta} \mathcal{V}_{r_{\hat{\omega}}}^{d_{\hat{\theta}}}(\rho)$ , where  $r_{\hat{\omega}}$  is the estimated reward. Then, with probability at least  $1 - \delta$ , the suboptimality gap incurred by RLHF is  $G(d_{\hat{\theta}}) = D(d_{\hat{\theta}}) + \tilde{\Theta}(\Lambda'_R \sqrt{d_R/n})$ .*

Now that we have a formulation for DPO, we can also extend the previous bounds for the MDP setting. Similar to Lemma 4.1, we now state an analogous result for this setting that guarantees the expression of the ground-truth reward in terms of optimal loglinear occupancy measures. Recall that  $\Phi$  and  $\Psi$  denote the reward and policy feature matrices, respectively. Let  $\pi_{r^*}^* \in \arg \max_{\pi} \mathcal{V}_{r^*}^{\pi}(\rho)$ . Define  $\Phi_{\pi_{r^*}^*}$  to be the  $d_R \times XY$ -dimensional matrix with columns defined as  $\gamma \mathbb{E}[\sum_t \gamma^t \phi(x_t, y_t) | x_0 = x, \pi_{r^*}^*] - \mathbb{E}[\sum_t \gamma^t \phi(x_t, y_t) | x_0 = x, y_0 = y, \pi_{r^*}^*]$ . We have the following result.

**Lemma 7.3.** *Assume that  $r^* \in \mathcal{F}$ ,  $d_{\rho}^{\mu} \in \Pi'$  and  $d_{\rho}^{\pi_{r^*}^*} \in \Pi'$ , for some optimal  $d_{\rho}^{\pi_{r^*}^*}$ . Furthermore, assume that the column space of  $\Phi + \Phi_{\pi_{r^*}^*}$  is contained in the column space of  $\Psi$ . Then, for finite MDPs with deterministic transitions, there exists  $\theta^*$  such that  $d_{\rho}^{\theta^*} \in \arg \max_d \mathcal{V}_{r^*}^d(\rho)$  and  $d_{\rho}^{\theta^*}(x, y) \propto d_{\rho}^{\mu}(x, y) \exp(A_{r^*}^{\pi_{r^*}^*}(x, y)/\beta)$ , where  $A_{r^*}^{\pi_{r^*}^*}(x, y) = Q_{r^*}^{\pi_{r^*}^*}(x, y) - V_{r^*}^{\pi_{r^*}^*}(x)$ .*

Now we are ready to state the DPO result for this section.

**Theorem 7.2.** *Let  $\delta > 0$ . Let  $d_{\rho}^{\tilde{\theta}}$  denote the occupancy measure returned by DPO and assume that  $d_{\rho}^{\pi_{r^*}^*} \in \Pi'$ ,*

*for some  $d_{\rho}^{\pi_{r^*}^*} \in \arg \max_{d_{\rho}} \mathcal{V}_{r^*}^{d_{\rho}}(\rho)$ , and that the condition of Lemma 7.3 is satisfied. Then, for any  $n \geq O(\text{tr}((\Sigma'_{\mathcal{D}_n, M})^\dagger)/(\beta(B')^2))$ , with probability at least  $1 - \delta$ , we have  $G(d_{\rho}^{\tilde{\theta}}) = D(d_{\rho}^{\tilde{\theta}}) + \Theta(\Lambda'_M(d_M + 1)/(\beta n) + \beta \Lambda_M \lambda (B')^2)$ .*

*Proof sketch.* We start by expressing the optimal discounted reward in terms of an optimal occupancy measure, which is also loglinear, by using Lemma 7.3. Then, we equivalently express the value function in terms of occupancy measures. This allows us to cancel out some terms and express the whole gap in terms of the  $D_{\text{KL}}(d_{\rho}^{\tilde{\theta}} || d_{\rho}^{\theta^*})$ . Finally, similar to the proof of Theorem 4.2, using loglinearity and properties of the log-sum-exp function, we obtain the desired bounds.  $\square$

### 7.4. Comparative Analysis

The main implication of the above results is that the observations made in Section 4 extend to deterministic MDPs, using our proposed formulation of RLHF and DPO. For the optimal value of  $\beta$  for DPO as discussed in Section 4.2, the RLHF and DPO bounds become directly comparable in terms of the dimension differences for deterministic MDPs. In MDPs with simple reward models (e.g., low-dimensional linear reward models), typically there is still a necessity for high-dimensional policy parameters to represent the value function effectively. This suggests that the complexity of the policy class exceeds that of the reward class and that RLHF outperforms DPO in such instances.

## 8. Concluding Discussion

In this paper, we provided a comparative analysis between reinforcement learning from human feedback (RLHF) and direct preference optimization (DPO). We performed a thorough analysis under different settings, where we derived sample complexity bounds for both paradigms and drew conclusions on their statistical comparison, based on sample size, regularization temperature, and the dimensionality of their respective parametrizations. We believe these results will initiate a larger discussion on the differences between these two paradigms.

There are many interesting future directions to pursue. The first natural extension of this work is to relax the assumptions made on policy and reward classes and provide a comparative analysis of RLHF and DPO on more realistic settings (eg. general function approximation). Next, a systematic large-scale empirical investigation that would validate the theoretical insights of this paper would be of great importance. Finally, as our current extension of DPO for MDPs is limited to deterministic MDPs using a loglinear occupancy measure regularization, it would be interesting to see whether DPO can be extended to more general formulations.

## Impact Statement

This paper focuses on the theoretical aspects of machine learning, providing a comparative analysis of different paradigms of learning from human preferences. We do not foresee any direct negative outcomes from the findings of this paper. On the contrary, we believe that our results might initiate a larger discussion on the statistical properties of learning from human preferences.

## Acknowledgements

The work of Andi Nika and Goran Radanovic was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – project number 467367360.

## References

- Agarwal, R., Schuurmans, D., and Norouzi, M. An Optimistic Perspective on Offline Reinforcement Learning. In *ICML*, 2020.
- Ailon, N., Karmin, Z. S., and Joachims, T. Reducing Dueling Bandits to Cardinal Bandits. In *ICML*, 2014.
- An, G., Lee, J., Zuo, X., Kosaka, N., Kim, K.-M., and Song, H. O. Direct Preference-based Policy Optimization without Reward Modeling. In *NeurIPS*, 2023.
- Azar, M. G., Rowland, M., Piot, B., Guo, D., Calandriello, D., Valko, M., and Munos, R. A General Theoretical Paradigm to Understand Learning from Human Preferences. *CoRR*, abs/2310.12036, 2023.
- Bai, Y. et al. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback. *CoRR*, abs/2204.05862, 2022.
- Bradley, R. A. and Terry, M. E. Rank Analysis of Incomplete Block Designs: I. The Method of Paired Comparisons. *Biometrika*, 39(3/4), 1952.
- Brown, D. S., Goo, W., Nagarajan, P., and Niekum, S. Extrapolating Beyond Suboptimal Demonstrations via Inverse Reinforcement Learning from Observations. In *ICML*, 2019.
- Chatterji, N., Pacchiano, A., Bartlett, P., and Jordan, M. On the Theory of Reinforcement Learning with Once-per-episode Feedback. In *NeurIPS*, 2021.
- Chen, X., Zhong, H., Yang, Z., Wang, Z., and Wang, L. Human-in-the-loop: Provably Efficient Preference-based Reinforcement Learning with General Function Approximation. In *ICML*, 2022.
- Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep Reinforcement Learning from Human Preferences. In *NeurIPS*, 2017.
- Faury, L., Abeille, M., Calauzènes, C., and Fercoq, O. Improved Optimistic Algorithms for Logistic Bandits. In *ICML*, 2020.
- Fujimoto, S., Meger, D., and Precup, D. Off-policy Deep Reinforcement Learning without Exploration. In *ICML*, 2019.
- Gajane, P., Urvoy, T., and Clérot, F. A Relative Exponential Weighing Algorithm for Adversarial Utility-based Dueling Bandits. In *ICML*, 2015.
- Ganguli, D. et al. Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned. *CoRR*, abs/2209.07858, 2022.
- Gao, L., Schulman, J., and Hilton, J. Scaling Laws for Reward Model Overoptimization. In *ICML*, 2023.
- Glaese, A. et al. Improving Alignment of Dialogue Agents via Targeted Human Judgements. *CoRR*, abs/2209.14375, 2022.
- Hejna, J., Rafailov, R., Sikchi, H., Finn, C., Niekum, S., Knox, W. B., and Sadigh, D. Contrastive Preference Learning: Learning from Human Feedback without RL. *CoRR*, abs/2310.13639, 2023.
- Jaques, N., Ghandeharioun, A., Shen, J. H., Ferguson, C., Lapedriza, A., Jones, N., Gu, S., and Picard, R. Way Off-policy Batch Deep Reinforcement Learning of Implicit Human Preferences in Dialog. *CoRR*, abs/1907.00456, 2019.
- Jin, Y., Yang, Z., and Wang, Z. Is Pessimism Provably Efficient for Offline RL? In *ICML*, 2021.
- Karimi, H., Nutini, J., and Schmidt, M. Linear Convergence of Gradient and Proximal-Gradient Methods under the Polyak-Łojasiewicz Condition. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I 16*. Springer, 2016.
- Kidambi, R., Rajeswaran, A., Netrapalli, P., and Joachims, T. Morel: Model-based Offline Reinforcement Learning. In *NeurIPS*, 2020.
- Knox, W. B. and Stone, P. Tamer: Training an Agent Manually via Evaluative Reinforcement. In *ICDL*, 2008.
- Komiyama, J., Honda, J., Kashima, H., and Nakagawa, H. Regret Lower Bound and Optimal Algorithm in Dueling Bandit Problem. In *COLT*, 2015.
- Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative Q-learning for Offline Reinforcement Learning. In *NeurIPS*, 2020.

- Laroche, R., Trichelair, P., and Des Combes, R. T. Safe Policy Improvement with Baseline Bootstrapping. In *ICML*, 2019.
- Lee, J., Jeon, W., Lee, B., Pineau, J., and Kim, K.-E. Optidice: Offline Policy Optimization via Stationary Distribution Correction Estimation. In *ICML*, 2021.
- MacGlashan, J., Ho, M. K., Loftin, R., Peng, B., Wang, G., Roberts, D. L., Taylor, M. E., and Littman, M. L. Interactive Learning from Policy-dependent Human Feedback. In *ICML*, 2017.
- Mei, J., Xiao, C., Szepesvári, C., and Schuurmans, D. On the Global Convergence Rates of Softmax Policy Gradient Methods. In *ICML*, 2020.
- Menick, J. et al. Teaching Language Models to Support Answers with Verified Quotes. *CoRR*, abs/2203.11147, 2022.
- Nachum, O., Dai, B., Kostrikov, I., Chow, Y., Li, L., and Schuurmans, D. Algaedice: Policy Gradient from Arbitrary Experience. *CoRR*, abs/1912.02074, 2019.
- Nakano, R. et al. Webgpt: Browser-assisted Question-answering with Human Feedback. *CoRR*, abs/2112.09332, 2021.
- Ouyang, L. et al. Training Language Models to Follow Instructions with Human Feedback. In *NeurIPS*, 2022.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., and Finn, C. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In *NeurIPS*, 2023.
- Ramamurthy, R., Ammanabrolu, P., Brantley, K., Hessel, J., Sifa, R., Bauckhage, C., Hajishirzi, H., and Choi, Y. Is Reinforcement Learning (not) for Natural Language Processing: Benchmarks, Baselines, and Building Blocks for Natural Language Policy Optimization. In *ICLR*, 2023.
- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. Bridging Offline Reinforcement Learning and Imitation Learning: A Tale of Pessimism. In *NeurIPS*, 2021.
- Saha, A. and Gopalan, A. Active Ranking with Subset-wise Preferences. In *AISTATS*, 2019.
- Saha, A. and Krishnamurthy, A. Efficient and Optimal Algorithms for Contextual Dueling Bandits under Realizability. In *ALT*, 2022.
- Saha, A., Pacchiano, A., and Lee, J. Dueling RL: Reinforcement Learning with Trajectory Preferences. In *AISTATS*, 2023.
- Shah, N. B., Balakrishnan, S., Bradley, J., Parekh, A., Ramch, K., Wainwright, M. J., et al. Estimation from Pairwise Comparisons: Sharp Minimax Bounds with Topology Dependence. In *AISTATS*, 2016.
- Shin, D., Dragan, A. D., and Brown, D. S. Benchmarks and Algorithms for Offline Preference-Based Reward Learning. *Transactions of Machine Learning Research*, 2023.
- Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. Learning to Summarize with Human Feedback. In *NeurIPS*, 2020.
- Uehara, M. and Sun, W. Pessimistic Model-based Offline Reinforcement Learning Under Partial Coverage. *CoRR*, abs/2107.06226, 2021.
- Wang, C. et al. Beyond Reverse KL: Generalizing Direct Preference Optimization with Diverse Divergence Constraints. *CoRR*, abs/2309.16240, 2023.
- Warnell, G., Waytowich, N., Lawhern, V., and Stone, P. Deep TAMER: Interactive Agent Shaping in High-dimensional State Spaces. In *AAAI*, 2018.
- Wu, J., Ouyang, L., Ziegler, D. M., Stiennon, N., Lowe, R., Leike, J., and Christiano, P. Recursively Summarizing Books with Human Feedback. *CoRR*, abs/2109.10862, 2021.
- Xie, T., Cheng, C.-A., Jiang, N., Mineiro, P., and Agarwal, A. Bellman-consistent Pessimism for Offline Reinforcement Learning. In *NeurIPS*, 2021.
- Xiong, W., Dong, H., Ye, C., Wang, Z., Zhong, H., Ji, H., Jiang, N., and Zhang, T. Iterative Preference Learning from Human Feedback: Bridging Theory and Practice for RLHF under KL-constraint. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, 2023.
- Yue, Y., Broder, J., Kleinberg, R., and Joachims, T. The K-armed Dueling Bandits Problem. In *COLT*, 2009.
- Zanette, A., Wainwright, M. J., and Brunskill, E. Provable Benefits of Actor-critic Methods for Offline Reinforcement Learning. In *NeurIPS*, 2021.
- Zhan, W., Uehara, M., Kallus, N., Lee, J. D., and Sun, W. Provable Offline Reinforcement Learning with Human Feedback. *CoRR*, abs/2305.14816, 2023.
- Zhu, B., Jordan, M. I., and Jiao, J. Principled Reinforcement Learning with Human Feedback from Pairwise or K-wise Comparisons. In *ICML*, 2023.

Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., Christiano, P., and Irving, G. Fine-tuning Language Models from Human Preferences. *CoRR*, abs/1909.08593, 2019.

Zoghi, M., Whiteson, S. A., De Rijke, M., and Munos, R. Relative Confidence Sampling for Efficient On-line Ranker Evaluation. In *WSDM*, 2014.

# Appendix

## Table of Contents

<b>A</b>	<b>Statistical Bounds for RLHF (Section 4 and 7)</b>	13
<b>B</b>	<b>Statistical Bounds for DPO (Section 4)</b>	16
<b>C</b>	<b>Statistical Bounds for DPO for MDPs (Section 7)</b>	17
<b>D</b>	<b>Convergence of Gradient Descent for RLHF Reward Learning (Section 5)</b>	20
<b>E</b>	<b>Convergence of Natural Policy Gradient for RLHF (Section 5)</b>	23
<b>F</b>	<b>Convergence of Gradient Descent for DPO (Section 5)</b>	29
<b>G</b>	<b>Non-realizable Rewards (Section 6)</b>	32
<b>H</b>	<b>The DPO Extension to MDPs</b>	33
<b>I</b>	<b>Gradient Expression for KL-regularized Objective in MDPs</b>	36
<b>J</b>	<b>Technical Lemmas</b>	37

### A. Statistical Bounds for RLHF (Section 4 and 7)

In this section, we prove the main RLHF result, Theorem 4.1. We state the detailed version of it together with the necessary constants.

**Theorem A.1.** *Let  $\delta > 0$ . Assume that the preference data satisfies the BT model, and  $r^* \in \mathcal{F}$ . Denote by  $\hat{\omega}$  and  $\hat{\theta}$  the reward and policy parameters learned via RLHF, respectively. Furthermore, assume that*

$$\hat{\omega} \in \arg \min_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$$

and

$$\hat{\theta} \in \arg \max_{\theta} \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\theta}}(\mathcal{D}_n).$$

Then, with probability at least  $1 - \delta$ , for any  $\lambda > 0$ , the suboptimality gap incurred by RLHF is

$$G(\pi_{\hat{\theta}}) = \Theta \left( \left\| (\Sigma_{\mathcal{D}_n, R} + \lambda I)^{-1/2} \right\|_2 \cdot \sqrt{\frac{d_R + \log(6/\delta)}{S_R^2 n} + \lambda F^2} \right) + D(\pi_{\hat{\theta}}),$$

where  $S_R = 1 / (2 + \exp(-2F) + \exp(2F))$ .

*Proof.* Let  $\Phi \in \mathbb{R}^{d_R \times XY}$  be the reward feature matrix. Then, for any  $\lambda > 0$ , with probability at least  $1 - \delta$ , we have

$$\begin{aligned} G(\pi_{\hat{\theta}}) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_{\hat{\theta}}}(\rho) \\ &= D(\pi_{\hat{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r^*}^{\pi_{\hat{\theta}}}(\rho) \right) \\ &= D(\pi_{\hat{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\hat{\theta}}}(\rho) \right) \\ &\stackrel{(a)}{\leq} D(\pi_{\hat{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\hat{\theta}}}(\rho) \right) \end{aligned}$$

$$\begin{aligned}
 &\stackrel{(b)}{\leq} D(\pi_{\hat{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r_{\hat{\omega}}}^*}(\mathcal{D}_n) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\mathcal{D}_n) \right) + O\left(\sqrt{\frac{\log(6/\delta)}{n}}\right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\hat{\theta}}}(\rho) \right) \\
 &\stackrel{(c)}{\leq} D(\pi_{\hat{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi_{r^*}^*}(\rho) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r^*}^*}(\rho) \right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\mathcal{D}_n) - \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{r_{\hat{\omega}}}^*}(\mathcal{D}_n) \right) + O\left(\sqrt{\frac{\log(6/\delta)}{n}}\right) + \left( \mathcal{V}_{r_{\hat{\omega}}}^{\pi_{\hat{\theta}}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\hat{\theta}}}(\rho) \right) \\
 &= D(\pi_{\hat{\theta}}) + \sum_{x,y} \rho(x) \cdot (\pi_{r^*}^*(y|x) - \pi_{\hat{\theta}}(y|x)) \cdot (r^*(x,y) - r_{\hat{\omega}}(x,y)) + O\left(\sqrt{\frac{\log(6/\delta)}{n}}\right) \\
 &\stackrel{(d)}{=} D(\pi_{\hat{\theta}}) + (d_{\rho}^* - d_{\rho}^{\pi_{\hat{\theta}}})^{\top} (\omega^* - \hat{\omega}) \Phi + O\left(\sqrt{\frac{\log(6/\delta)}{n}}\right) \\
 &\stackrel{(e)}{\leq} D(\pi_{\hat{\theta}}) + \|\Phi (d_{\rho}^* - d_{\rho}^{\pi_{\hat{\theta}}})\|_{(\Sigma_{\mathcal{D}_n, R} + \lambda I)^{-1}} \|\omega^* - \hat{\omega}\|_{\Sigma_{\mathcal{D}_n, R} + \lambda I} + O\left(\sqrt{\frac{\log(6/\delta)}{n}}\right) \\
 &\stackrel{(f)}{\leq} D(\pi_{\hat{\theta}}) + O\left(\left\|(\Sigma_{\mathcal{D}_n, R} + \lambda I)^{-1/2}\right\|_2 \cdot \sqrt{\frac{d_R + \log(6/\delta)}{S_R^2 n}} + \lambda F^2\right),
 \end{aligned}$$

where (a) is due to the fact that  $\pi_{r_{\hat{\omega}}}^* \in \arg \max_{\pi} \mathcal{V}_{r_{\hat{\omega}}}^{\pi}(\rho)$ ; (b) follows from Lemma A.2 and the union bound – if we have  $\mathbb{P}(\mathcal{E}_i^c) \leq \delta_i$ , for  $i = 1, 2, 3$ , where  $\mathcal{E}^c$  denotes the complement of event  $\mathcal{E}$ , letting  $\delta_i = \delta/3$ , for all  $i$ , we have

$$\begin{aligned}
 \mathbb{P}(\mathcal{E}_1 \cup \mathcal{E}_2 \cup \mathcal{E}_3) &= 1 - \mathbb{P}(\mathcal{E}_1^c \cup \mathcal{E}_2^c \cup \mathcal{E}_3^c) \\
 &\geq 1 - (\mathbb{P}(\mathcal{E}_1) + \mathbb{P}(\mathcal{E}_2) + \mathbb{P}(\mathcal{E}_3)) \\
 &\geq 1 - (\delta_1 + \delta_2 + \delta_3) \\
 &= 1 - \delta.
 \end{aligned}$$

Next, (c) is due to the fact that  $\pi_{\hat{\theta}} \in \arg \max_{\pi \in \Pi} \mathcal{V}_{r_{\hat{\omega}}}^{\pi}(\mathcal{D}_n)$  and  $\pi_{r_{\hat{\omega}}}^* \in \Pi$  (as per Lemma 4.1); (d) is due to  $d_{\rho}^*(x, y) = \rho(x, y)\pi_{r^*}^*(y|x)$  and  $d_{\rho}^{\pi_{\hat{\theta}}}(x, y) = \rho(x)\pi_{\hat{\theta}}(y|x)$ ; (e) is an application of the Cauchy-Schwarz inequality with respect to the semi-norm induced by matrix  $\Sigma_{\mathcal{D}_n, R} + \lambda I$ ; and (f) is a direct application of Lemma 3.1 of (Zhu et al., 2023) for the discounted infinite-horizon setting.

The lower bound is an immediate application of Theorem 3.10 of (Zhu et al., 2023). Note that we are under the same conditions; our reward function is assumed to be linear, and we also assume a bounded covering number. For the lower bound construction, let  $\mu = \pi_{r^*}^{\text{opt}}$ , i.e., the reference policy is the actually optimal one. Let  $\text{CB}(\Lambda)$  denote the set of bandit instances coupled with datasets with a covering number no more than  $\Lambda$ . Let  $\mathcal{Q}$  denote such an instance. Under these assumptions, Theorem 3.10 of (Zhu et al., 2023) implies an information-theoretic lower bound of

$$\inf_{\pi} \sup_{\mathcal{Q} \in \text{CB}(\Lambda_R)} (V_{\mathcal{Q}}^{\text{opt}}(\rho) - V_{\mathcal{Q}}^{\pi}(\rho)) \geq O\left(\Lambda_R \sqrt{\frac{d_R}{n}}\right).$$

□

**Theorem 7.1.** *Let  $\delta > 0$ . Assume that the policy learning phase yields  $\hat{\theta} \in \arg \max_{\theta} \mathcal{V}_{r_{\hat{\omega}}}^{d_{\rho}^{\theta}}(\rho)$ , where  $r_{\hat{\omega}}$  is the estimated reward. Then, with probability at least  $1 - \delta$ , the suboptimality gap incurred by RLHF is  $G(d_{\rho}^{\hat{\theta}}) = D(d_{\rho}^{\hat{\theta}}) + \tilde{\Theta}(\Lambda'_R \sqrt{d_R/n})$ .*

*Proof.* The proof of this result is an immediate application of the previous result with instead an application of Lemma 5.1 of (Zhu et al., 2023) with

$$S_R = 1 / (2 + \exp(-2F(1 - \gamma)) + \exp(2F(1 - \gamma))).$$

□

Next, we connect the difference between the sample regularized gap and the expected regularized gap with respect to the context distribution, and obtain the following result.

**Lemma A.2.** Let  $\delta > 0$  and assume that the conditions of Theorem A.1 are satisfied. Then, we have that

$$\left| \mathbb{E}_{x \sim \rho} \left[ \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{r_{\widehat{\omega}}}^*}(x) \right] - \frac{1}{n} \sum_{x \in \mathcal{D}_n} \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{r_{\widehat{\omega}}}^*}(x) \right| \leq \sqrt{\frac{\log(4/\delta)}{n}},$$

and

$$\left| \frac{1}{n} \sum_{x \in \mathcal{D}_n} \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{\widehat{\theta}}}(x) - \mathbb{E}_{x \sim \rho} \left[ \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{\widehat{\theta}}}(x) \right] \right| \leq \sqrt{\frac{(1 + \beta(2B + \log Y)) \log(4/\delta)}{n}}.$$

with probability at least  $1 - \delta$ .

*Proof.* Using the reward-to-policy mapping of Equation (1), we have that, for every  $(x, y)$ ,

$$\mu(y|x) = \pi_{r_{\widehat{\omega}}}^*(y|x) \exp\left(-\frac{1}{\beta} r_{\widehat{\omega}}(x, y)\right).$$

Thus, note that, for every  $x \in \mathcal{X}$ ,

$$\begin{aligned} \left| \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{r_{\widehat{\omega}}}^*}(x) \right| &= \left| \sum_y \pi_{r_{\widehat{\omega}}}^*(y|x) r_{\widehat{\omega}}(x, y) - \beta \pi_{r_{\widehat{\omega}}}^*(y|x) \log \frac{\pi_{r_{\widehat{\omega}}}^*(y|x)}{\mu(y|x)} \right| \\ &\leq 1 + \beta \left| \sum_y \pi_{r_{\widehat{\omega}}}^*(y|x) \log \frac{\pi_{r_{\widehat{\omega}}}^*(y|x)}{\mu(y|x)} \right| \\ &= 1 + \beta \left| \sum_y \pi_{r_{\widehat{\omega}}}^*(y|x) \frac{1}{\beta} r_{\widehat{\omega}}(x, y) \right| \\ &\leq 2, \end{aligned}$$

where we have used that the reward lies in  $[0, 1]$ . On the other hand, we have

$$\begin{aligned} \left| \mathcal{V}_{r_{\widehat{\omega}}}^{\pi_{\widehat{\theta}}}(x) \right| &= \left| \sum_y \pi_{\widehat{\theta}}(y|x) r_{\widehat{\omega}}(x, y) - \beta \pi_{\widehat{\theta}}(y|x) \log \frac{\pi_{\widehat{\theta}}(y|x)}{\mu(y|x)} \right| \\ &\leq 1 + \beta \left| \sum_y \pi_{\widehat{\theta}}(y|x) \left( \log \frac{\pi_{\widehat{\theta}}(y|x)}{\pi_{r_{\widehat{\omega}}}^*(y|x)} + \frac{1}{\beta} r_{\widehat{\omega}}(x, y) \right) \right| \\ &\leq 2 + \beta \max_y \left| \log \frac{\pi_{\widehat{\theta}}(y|x)}{\pi_{\theta^*}(y|x)} \right| \\ &\leq 2 + \beta \max_{x, y} \left( \left| \log \frac{\exp(\psi(x, y)^\top \widehat{\theta})}{\sum_{y'} \exp(\psi(x, y')^\top \widehat{\theta})} \right| + \left| \log \frac{\exp(\psi(x, y)^\top \theta^*)}{\sum_{y'} \exp(\psi(x, y')^\top \theta^*)} \right| \right) \\ &\leq 2 + \beta \max_{x, y} \left( \left| \log \exp(\psi(x, y)^\top \widehat{\theta}) \right| + \left| \log \sum_{y'} \exp(\psi(x, y')^\top \widehat{\theta}) \right| \right. \\ &\quad \left. + \left| \log \exp(\psi(x, y)^\top \theta^*) \right| + \left| \log \sum_{y'} \exp(\psi(x, y')^\top \theta^*) \right| \right) \\ &\leq 2 + \beta (2B + 2 \log(Y \exp(B))) \\ &\leq 2 + 2\beta(2B + \log Y), \end{aligned}$$

where we have used Lemma J.1 and the fact that

$$-B \leq \langle \psi(x, y), \theta \rangle \leq B.$$

The result then follows from Hoeffding's inequality.  $\square$

## B. Statistical Bounds for DPO (Section 4)

In this section, we prove the main DPO result for Section 4.

**Theorem 4.2.** *Let  $\delta > 0$  and  $\beta > 0$ . Assume that  $r^* \in \mathcal{F}$ ,  $\mu \in \Pi$ , and that the condition of Lemma 4.1 is satisfied. Let  $n \geq O\left(\text{tr}(\Sigma_{\mathcal{D}_{n,P}}^\dagger)/(\beta B^2)\right)$ . Then, with probability at least  $1 - \delta$ , the suboptimality gap of DPO is*

$$G(\pi_{\tilde{\theta}}) = D(\pi_{\tilde{\theta}}) + \Theta\left(\frac{\Lambda_P(d_P + 1)}{\beta n} + \beta \lambda \Lambda_P B^2\right).$$

*Proof.* Lemma 4.1 implies that there exists  $\theta^* \in \mathbb{R}^d$ , such that, for every  $(x, y)$ ,

$$r^*(x, y) = \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} + \beta Z(x)$$

and  $\mathcal{V}_{r^*}(\rho) = \mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho)$ . Now, observe that

$$\begin{aligned} G(\pi_{\tilde{\theta}}) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_{\tilde{\theta}}}(\rho) \\ &= D(\pi_{\tilde{\theta}}) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\tilde{\theta}}}(\rho)) \\ &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\theta^*}(\cdot|x)}} \left[ r^*(x, y) - \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} \right] - \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\tilde{\theta}}(\cdot|x)}} \left[ r^*(x, y) - \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} \right] \\ &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\theta^*}(\cdot|x)}} \left[ \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} + \beta \log Z(x) - \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} \right] \\ &\quad - \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\tilde{\theta}}(\cdot|x)}} \left[ \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} + \beta \log Z(x) - \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} \right] \\ &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\tilde{\theta}}(\cdot|x)}} \left[ \beta \log \pi_{\tilde{\theta}}(y|x) - \beta \log \pi_{\theta^*}(y|x) \right] \\ &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\tilde{\theta}}(\cdot|x)}} \left[ \beta \langle \psi(x, y), \tilde{\theta} - \theta^* \rangle + \beta \log \frac{\sum_{y'} \exp(\psi(x, y')^\top \theta^*)}{\sum_{y'} \exp(\psi(x, y')^\top \tilde{\theta})} \right] \\ &= D(\pi_{\tilde{\theta}}) + \beta \mathbb{E}_{\substack{x \sim \rho \\ y \sim \pi_{\tilde{\theta}}(\cdot|x)}} \left[ \langle \psi(x, y), \tilde{\theta} - \theta^* \rangle \right] + \beta (A(\theta^*) - A(\tilde{\theta})), \end{aligned}$$

where we have denoted by  $A(\theta)$  the log-sum-exp function

$$A(\theta) = \sum_x \rho(x) \log \sum_{y'} \exp(\psi(x, y')^\top \theta),$$

At this point, some properties of the log-exp-sum function will be useful. The proof of the following result can be found in Appendix J.

**Lemma B.1.** *The function  $A(\theta)$  is 1-Lipschitz and 2-smooth. Moreover, if the features are sampled from a 0-mean distribution and span  $R^{d_P}$ , then there exists  $\kappa > 0$ , such that  $A(\theta)$  is  $\kappa$ -strongly convex.*

Since  $A(\theta)$  is 2-smooth, we have

$$\begin{aligned} A(\theta^*) - A(\tilde{\theta}) &\leq \langle \nabla_\theta A(\tilde{\theta}), \theta^* - \tilde{\theta} \rangle + \|\theta^* - \tilde{\theta}\|_2^2 \\ &= \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \langle \psi(x, y), \theta^* - \tilde{\theta} \rangle \right] + \|\theta^* - \tilde{\theta}\|_2^2. \end{aligned}$$

Substituting to the suboptimality gap equalities, we obtain

$$\begin{aligned} G(\pi_{\tilde{\theta}}) &\leq D(\pi_{\tilde{\theta}}) + \|\theta^* - \tilde{\theta}\|_2^2 \\ &\stackrel{(a)}{\leq} D(\pi_{\tilde{\theta}}) + \beta \left\| (\Sigma_{\mathcal{D}_{n,P}} + \lambda I)^{-1} \right\|_2 \|\theta^* - \tilde{\theta}\|_{\Sigma_{\mathcal{D}_{n,P}} + \lambda I}^2 \end{aligned}$$



$$\begin{aligned}
 &\stackrel{(b)}{\leq} D(\pi_{\tilde{\theta}}) + \beta \left\| (\Sigma_{\mathcal{D}_n, P} + \lambda I)^{-1} \right\|_2 \left\| \theta^* - \tilde{\theta} \right\|_{\Sigma_{\mathcal{D}_n, P}}^2 + 4\beta\lambda\Lambda_P B^2 \\
 &\stackrel{(c)}{\leq} D(\pi_{\tilde{\theta}}) + O\left(\frac{\Lambda_P(d_P + 1)}{\beta n}\right) + 4\beta\lambda\Lambda_P B^2,
 \end{aligned}$$

where for (a) we have used that  $\langle x, Ax \rangle \leq \|A\|_2 \|x\|_2^2$ ; (b) is due to the fact that  $\|\theta\|_2 \leq B$ ; (c) follows from Theorem J.7 and the assumption that  $\tilde{\theta} \in \arg \min_{\theta} \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)$ . The value of  $\lambda$  can be tuned accordingly. Note that, for fixed  $\beta$ , letting  $\lambda = \Theta(1/n)$  yields the desired bound. If  $\beta = \Theta(1/\sqrt{n})$ , then any small value of  $\lambda$  works.

On the other hand, consider the following feature construction. Let  $\Psi$  be full rank with zero-mean columns. Then, there exists  $\kappa > 0$  such that  $A(\theta)$  is  $\kappa$ -strongly convex. This, in turn, implies that

$$\begin{aligned}
 A(\theta^*) - A(\tilde{\theta}) &\geq \left\langle \nabla_{\theta} A(\tilde{\theta}), \theta^* - \tilde{\theta} \right\rangle + \frac{\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_2^2 \\
 &= -\mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \left\langle \psi(x, y), \tilde{\theta} - \theta^* \right\rangle \right] + \frac{\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_2^2.
 \end{aligned}$$

Thus, substituting in the original gap expression, we obtain

$$\begin{aligned}
 G(\pi_{\tilde{\theta}}) &\geq D(\pi_{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_2^2 \\
 &\geq D(\pi_{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_2^2 \|\Sigma_{\mathcal{D}_n, P}\|_2 \\
 &\geq D(\pi_{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\langle \tilde{\theta} - \theta^*, \Sigma_{\mathcal{D}_n, P}(\tilde{\theta} - \theta^*) \right\rangle \\
 &\geq D(\pi_{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_{\Sigma_{\mathcal{D}_n, P}} \\
 &\geq D(\pi_{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_{\Sigma_{\mathcal{D}_n, P}} \\
 &\geq D(\pi_{\tilde{\theta}}) + \Omega\left(\frac{(d_P + 1)}{\beta n}\right),
 \end{aligned}$$

for any  $n \geq O\left(\frac{\text{tr}(\Sigma_{\mathcal{D}_n, P}^{\dagger})}{(\beta B^2)}\right)$ , where the second inequality uses the fact that  $\|\Sigma_{\mathcal{D}_n, P}\|_2 \leq 1$ ; the third inequality follows from Cauchy-Schwarz; the last inequality follows by Theorem J.7.  $\square$

### C. Statistical Bounds for DPO for MDPs (Section 7)

In this section, we will prove Corollary 7.2, the statistical convergence rate of the DPO method in the MDP setting. We restate the result with all the quantities appearing in the bounds.

**Theorem C.1.** *Assume that  $r^* \in \mathcal{F}$ , the data  $\mathcal{D}$  satisfies the BT preference model for trajectories, and that the feature matrix is full rank. Furthermore, assume that  $d_{\rho}^* \in \Pi'$  and assume that we have 0 optimization error from gradient descent on the MLE loss. Let  $\beta > 0$ , and*

$$\begin{aligned}
 S_M &= (\exp(-B') + \exp(B') + 2)^{-1}, \\
 U' &= \exp(-2B') + \exp(2B') + 2, \\
 \Lambda_M &= \left\| (\Sigma'_{\mathcal{D}_n, M} + \lambda I)^{-1/2} \right\|_2.
 \end{aligned}$$

Then, DPO incurs the following minimax bounds on the suboptimality gap:

$$G(d_{\rho}^{\tilde{\theta}}) = D(d_{\rho}^{\tilde{\theta}}) + \Theta\left(\frac{\Lambda_M U' (d_M + 1)}{\beta S_P n}\right).$$

*Proof.* First, note that, under some assumptions on the feature space, Lemma J.2 implies that there exists  $\theta^*$  for which we

can have  $\mathcal{V}_{r^*}^*(\rho) = \mathcal{V}_{r^*}^{d_{\rho}^{\theta^*}}(\rho)$ , with respect to some policy  $\pi_{r^*}$  and that

$$\sum_{t \geq 0} \gamma^t r^*(x_t, y_t) = \sum_{t \geq 0} \gamma^t \beta \log \frac{d_{\rho}^{\theta^*}(x_t, y_t)}{d_{\rho}^{\mu}(x_t, y_t)} + \beta + \alpha^*(x_0), \quad (7)$$

for any trajectory  $\tau$ , where  $\alpha^*$  is the optimal dual variable of Problem (P3.2'), and we have used Equation (15) to express the ground-truth reward in terms of an optimal regularized policy. Now, let us denote by  $\tilde{\pi}$  the policy corresponding to  $d_{\rho}^{\tilde{\theta}}$ . Observe that

$$\begin{aligned} G(d_{\rho}^{\tilde{\theta}}) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{d_{\rho}^{\tilde{\theta}}}(\rho) \\ &= D(d_{\rho}^{\tilde{\theta}}) + \left( \mathcal{V}_{r^*}^{d_{\rho}^{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{d_{\rho}^{\tilde{\theta}}}(\rho) \right) \\ &= D(d_{\rho}^{\tilde{\theta}}) + \mathbb{E}_{x_0 \sim \rho, y_t \sim \pi_{r^*}(\cdot | x_t)} \left[ \sum_{t \geq 0} \gamma^t r^*(x_t, y_t) \right] - \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} \left[ \log \frac{d_{\rho}^{\theta^*}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] - \mathcal{V}_{r^*}^{d_{\rho}^{\tilde{\theta}}}(\rho) \\ &= D(d_{\rho}^{\tilde{\theta}}) + \mathbb{E}_{x_0 \sim \rho, y_t \sim \pi_{r^*}(\cdot | x_t)} \left[ \sum_{t \geq 0} \gamma^t \beta \log \frac{d_{\rho}^{\theta^*}(x_t, y_t)}{d_{\rho}^{\mu}(x_t, y_t)} + \beta + \alpha^*(x_0) \right] - \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} \left[ \log \frac{d_{\rho}^{\theta^*}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] - \mathcal{V}_{r^*}^{d_{\rho}^{\tilde{\theta}}}(\rho) \end{aligned} \quad (8)$$

$$= D(d_{\rho}^{\tilde{\theta}}) + \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} \left[ \beta \log \frac{d_{\rho}^{\theta^*}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] + \frac{\beta}{1-\gamma} + \sum_x \rho(x) \alpha^*(x) - \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} \left[ \log \frac{d_{\rho}^{\theta^*}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] - \mathcal{V}_{r^*}^{d_{\rho}^{\tilde{\theta}}}(\rho) \quad (9)$$

$$\begin{aligned} &= D(d_{\rho}^{\tilde{\theta}}) + \frac{\beta}{1-\gamma} \\ &\quad + \sum_x \rho(x) \alpha^*(x) - \mathbb{E}_{x_0 \sim \rho, y_t \sim \tilde{\pi}(\cdot | x_t)} \left[ \sum_{t \geq 0} \gamma^t \beta \log \frac{d_{\rho}^{\theta^*}(x_t, y_t)}{d_{\rho}^{\mu}(x_t, y_t)} + \beta + \alpha^*(x_0) \right] + \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \log \frac{d_{\rho}^{\tilde{\theta}}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] \\ &= D(d_{\rho}^{\tilde{\theta}}) + \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \log \frac{d_{\rho}^{\tilde{\theta}}(x, y)}{d_{\rho}^{\mu}(x, y)} - \log \frac{d_{\rho}^{\theta^*}(x, y)}{d_{\rho}^{\mu}(x, y)} \right] \\ &= D(d_{\rho}^{\tilde{\theta}}) + \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \beta \log d_{\rho}^{\tilde{\theta}}(x, y) - \beta \log d_{\rho}^{\theta^*}(x, y) \right] \\ &= D(d_{\rho}^{\tilde{\theta}}) + \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \beta \left\langle \psi'(x, y), \tilde{\theta} - \theta^* \right\rangle + \beta \log \frac{\sum_{x', y'} \exp(\psi'(x', y')^\top \theta^*)}{\sum_{x', y'} \exp(\psi'(x', y')^\top \tilde{\theta})} \right] \\ &= D(d_{\rho}^{\tilde{\theta}}) + \beta \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \left\langle \psi'(x, y), \tilde{\theta} - \theta^* \right\rangle + A(\theta^*) - A(\tilde{\theta}) \right], \end{aligned} \quad (10)$$

where we have denoted by

$$A(\theta) = \log \sum_{x', y'} \exp(\psi'(x', y')^\top \theta)$$

the log-sum-exp function. Above, in Equation (8) we have applied Equation (7); Equation (9) uses the fact that, for any policy  $\pi$  and function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , we have

$$\mathbb{E}_{x \sim \rho, y_t \sim \pi(\cdot | x_t)} \left[ \sum_{t \geq 0} \gamma^t f(x_t, y_t) \right] = \mathbb{E}_{(x,y) \sim d^{\pi}} [f(x, y)].$$

Finally, for Equation (10) we have used loglinearity. Now, given  $\theta \in \mathbb{R}^{d_P}$ , note that

$$\nabla_{\theta} A(\theta) = \frac{\sum_{x,y} \exp(\psi'(x, y)^\top \theta) \cdot \psi'(x, y)}{\sum_{x', y'} \exp(\psi'(x', y')^\top \theta)} = \sum_{x,y} d_{\rho}^{\theta}(x, y) \psi'(x, y).$$

On the other hand, the Hessian of  $A(\theta)$  is

$$\begin{aligned}
 \nabla_{\theta}^2 A(\theta) &= \sum_{x,y} \nabla_{\theta} d_{\rho}^{\theta}(x,y) \psi'(x,y) \\
 &= \sum_{x,y} d_{\rho}^{\theta}(x,y) \left( \psi'(x,y) - \mathbb{E}_{(x',y') \sim d_{\rho}^{\theta}} [\psi'(x',y')] \right) \psi'(x,y)^{\top} \\
 &= \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta}} \left[ \psi'(x,y) \psi'(x,y)^{\top} \right] - \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta}} [\psi'(x,y)] \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta}} [\psi'(x,y)]^{\top} \\
 &= \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta}} \left[ \left( \psi'(x,y) - \mathbb{E}_{\theta} [\psi'(x,y)] \right) \left( \psi'(x,y) - \mathbb{E}_{\theta} [\psi'(x,y)] \right)^{\top} \right].
 \end{aligned}$$

By assumption on the feature mapping, we have that

$$\begin{aligned}
 \|\nabla_{\theta}^2 A(\theta)\|_2 &\leq \max_{x,y} \left\| \left( \psi'(x,y) - \mathbb{E}_{\theta} [\psi'(x,y)] \right) \left( \psi'(x,y) - \mathbb{E}_{\theta} [\psi'(x,y)] \right)^{\top} \right\|_2 \\
 &\leq \max_{x,y} \|\psi'(x,y) - \mathbb{E}_{\theta} [\psi'(x,y)]\|_2 \\
 &\leq 2 \max_{x,y} \|\psi'(x,y)\|_2 = 2.
 \end{aligned}$$

Therefore, the function  $A(\theta)$  is 2-smooth in  $\theta$ , which implies that

$$\begin{aligned}
 A(\theta^*) - A(\tilde{\theta}) &\leq \left\langle \nabla_{\theta} A(\tilde{\theta}), \theta^* - \tilde{\theta} \right\rangle + \|\theta^* - \tilde{\theta}\|_2^2 \\
 &= \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} \left[ \left\langle \psi'(x,y), \theta^* - \tilde{\theta} \right\rangle \right] + \|\theta^* - \tilde{\theta}\|_2^2.
 \end{aligned}$$

Substituting to the suboptimality gap equalities, we obtain

$$\begin{aligned}
 G(d_{\rho}^{\tilde{\theta}}) &\leq D(d_{\rho}^{\tilde{\theta}}) + \|\theta^* - \tilde{\theta}\|_2^2 \\
 &\leq D(d_{\rho}^{\tilde{\theta}}) + \beta \left\| (\Sigma'_{\mathcal{D}_n, M} + \lambda I)^{-1} \right\|_2 \|\theta^* - \tilde{\theta}\|_{\Sigma'_{\mathcal{D}_n, M} + \lambda I}^2 \\
 &\leq D(d_{\rho}^{\tilde{\theta}}) + \frac{\Lambda_M U' d_M}{S_M \beta n} + 4\beta \Lambda_M \lambda (B')^2,
 \end{aligned}$$

where the last inequality follows from Theorem J.7 and the assumption on exact optimization.

For the lower bound, let  $\psi'$  be sampled from a 0-mean bounded distribution. Note that, for any non-zero vector in  $\mathbb{R}^{d_M}$ , we have

$$\begin{aligned}
 z^{\top} \nabla_{\theta}^2 A(\theta) z &= \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta}} \left[ z^{\top} \psi'(x,y) \psi'(x,y)^{\top} z \right] \\
 &\geq \min_{\theta, x,y} d_{\rho}^{\theta}(x,y) \sum_{x,y} (\psi'(x,y)^{\top} z)^2 \\
 &\geq C_3 \sum_{x,y} (\psi'(x,y)^{\top} z)^2,
 \end{aligned}$$

for a positive  $C_3$ , since  $d_{\rho}^{\theta}$  is in the loglinear class, for every  $\theta$ . Now, note that, if  $z$  can be expressed as a linear combination of  $\{\psi'(x,y)\}_{x,y}$ , the summation cannot be zero for non-zero  $z$ . Thus, if  $\{\psi'(x,y)\}_{x,y}$  spans  $\mathbb{R}^{d_M}$ , that is, the feature matrix is full rank, then there exists an absolute positive constant  $\kappa$ , such that we have

$$\|\nabla_{\theta}^2 A(\theta)\|_2 \geq \kappa > 0.$$

Thus, the function  $A(\theta)$  is  $\kappa$ -strongly convex. This, in turn, implies that

$$\begin{aligned}
 A(\theta^*) - A(\tilde{\theta}) &\geq \left\langle \nabla_{\theta} A(\tilde{\theta}), \theta^* - \tilde{\theta} \right\rangle + \frac{\kappa}{2} \|\tilde{\theta} - \theta^*\|_2^2 \\
 &\geq \left\langle \nabla_{\theta} A(\tilde{\theta}), \theta^* - \tilde{\theta} \right\rangle + \frac{\kappa}{2} \|\tilde{\theta} - \theta^*\|_{\Sigma'_{\mathcal{D}_n, M}}^2
 \end{aligned}$$

$$= -\mathbb{E}_{(x,y) \sim d_{\tilde{\rho}}} \left[ \left\langle \psi'(x,y), \tilde{\theta} - \theta^* \right\rangle \right] + \frac{\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_{\Sigma'_{\mathcal{D}_n, M}}^2,$$

using a similar argument as in the proof of Theorem 4.2. Thus, substituting in the original gap expression, we obtain

$$\begin{aligned} G(d_{\tilde{\rho}}^{\tilde{\theta}}) &\geq D(d_{\tilde{\rho}}^{\tilde{\theta}}) + \frac{\beta\kappa}{2} \left\| \tilde{\theta} - \theta^* \right\|_2^2 \\ &\geq D(d_{\tilde{\rho}}^{\tilde{\theta}}) + \Omega\left(\frac{d_M + 1}{\beta n}\right), \end{aligned}$$

for any  $n \geq O(\text{tr}((\Sigma'_D)^\dagger)/(S_M(B')^2))$ , by Theorem J.7.  $\square$

## D. Convergence of Gradient Descent for RLHF Reward Learning (Section 5)

In this section, we will prove convergence bounds for projected gradient descent for the RLHF reward learning phase. Recall that, the projected gradient update rule for the reward learning phase is given as

$$\omega_{t+1} = \text{proj}_{\omega: \|\omega\|_2 \leq F} (\omega_t - \eta \nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)).$$

First, we show Lipschitzness and smoothness of the loss function.

**Lemma D.1.** *The RLHF reward learning objective  $\mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$  is  $2 \exp(2F)$ -Lipschitz and  $2 \exp(2F)$ -smooth.*

*Proof.* Note that the gradient of  $\mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$  satisfies

$$\begin{aligned} \|\nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)\|_2 &= \left\| \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \nabla_{\omega} \log(1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))) \right\|_2 \\ &\leq \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \frac{\exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))}{1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))} \|\phi(x, y^w) - \phi(x, y^l)\|_2 \\ &\leq \exp(2F) \|\phi(x, y^w) - \phi(x, y^l)\|_2 \\ &\leq 2 \exp(2F). \end{aligned}$$

Moreover, the Hessian of  $\mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$  satisfies

$$\begin{aligned} \|\nabla_{\omega}^2 \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)\|_2 &= \left\| \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \nabla_{\omega} \frac{\exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))}{1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))} (\phi(x, y^w) - \phi(x, y^l)) \right\|_2 \\ &\leq \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \frac{\exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))}{(1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l))))^2} \left\| (\phi(x, y^w) - \phi(x, y^l)) (\phi(x, y^w) - \phi(x, y^l))^\top \right\|_2 \\ &\leq \exp(2F) \|\phi(x, y^w) - \phi(x, y^l)\|_2 \\ &\leq 2 \exp(2F). \end{aligned}$$

The result follows.  $\square$

Next, we show that  $\mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$  satisfies the PL condition (Karimi et al., 2016) defined below.

**Definition D.2.** A function  $\mathcal{L}$  is said to satisfy the PL condition with coefficient  $C_{PL} > 0$  if, for every  $\omega$  in the domain of  $\mathcal{L}$ , we have

$$\|\nabla_{\omega} \mathcal{L}(\omega)\|_2^2 \geq C_{PL} (\mathcal{L}(\omega) - \mathcal{L}^*),$$

where  $\mathcal{L}^*$  is the minimum value of  $\mathcal{L}$  in its domain.

**Lemma D.3.** Let  $L_2 = 2 \exp(2F)$  and

$$C_{PL} = \frac{\exp(-2F)\xi(1 + \exp(-2F))}{n(1 + \exp(2F))^2},$$

where

$$0 < \xi = \min_{(x, y^w, y^l) \sim \mathcal{D}_n} \|\phi(x, y^w) - \phi(x, y^l)\|_2^2.$$

Then, we have

$$\frac{1}{2} \|\nabla_\omega \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)\|^2 \geq C_{PL} (\mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n) - \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n)).$$

*Proof.* Due to the assumption on the features and parameter vectors, we have

$$\begin{aligned} \frac{1}{2} \|\nabla_\omega \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n)\|^2 &\geq \frac{1}{2n^2} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \frac{\exp(-2F)}{1 + \exp(2F)} \min_{(x, y^w, y^l) \sim \mathcal{D}_n} \|\phi(x, y^w) - \phi(x, y^l)\|_2^2 \\ &\geq \frac{\exp(-2F)\xi}{n(1 + \exp(2F))}. \end{aligned}$$

On the other hand, note that, for some  $\omega^*$  such that  $\|\omega^*\|_2 \leq F$ , we have

$$\begin{aligned} \mathcal{L}_{\text{RLHF}}^\omega(\mathcal{D}_n) - \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n) &= \mathbb{E}_{(x, y^w, y^l) \sim \mathcal{D}_n} \left[ \log(1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))) \right. \\ &\quad \left. - \log(1 + \exp((\omega^*)^\top (\phi(x, y^w) - \phi(x, y^l)))) \right] \\ &= \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \log \frac{1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))}{1 + \exp((\omega^*)^\top (\phi(x, y^w) - \phi(x, y^l)))} \\ &\leq \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \left( \frac{1 + \exp(\omega^\top (\phi(x, y^w) - \phi(x, y^l)))}{1 + \exp((\omega^*)^\top (\phi(x, y^w) - \phi(x, y^l)))} - 1 \right) \\ &\leq \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \frac{1 + \exp(2F)}{1 + \exp(-2F)} \\ &\leq \frac{1 + \exp(2F)}{1 + \exp(-2F)}, \end{aligned}$$

where the third inequality follows from  $\log x \leq x - 1$ , for  $x > 0$ . Solving for  $C_{PL}$  the equation

$$C_{PL} \frac{1 + \exp(2F)}{1 + \exp(-2F)} = \frac{\exp(-2F)\xi}{n(1 + \exp(2F))},$$

we obtain

$$C_{PL} = \frac{\exp(-2F)\xi(1 + \exp(-2F))}{n(1 + \exp(2F))^2}.$$

□

Now, we are ready to state the convergence result for gradient descent.

**Theorem 5.1.** For every  $t \geq 0$ , the gradient descent procedure (4) with learning rate  $\eta = 1/\exp(2F)$  satisfies

$$\|\omega_t - \omega_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, R}}^2 \leq O\left(1 - \frac{1}{n}\right)^t.$$

*Proof.* The projected gradient descent rule is equivalent to the proximal gradient update (Karimi et al., 2016), given as

$$\omega_{t+1} = \arg \min_{\omega} \left( \langle \nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega_t}(\mathcal{D}_n), \omega - \omega_t \rangle + \frac{L_2}{2} \|\omega - \omega_t\|_2^2 + g(\omega) - g(\omega_t) \right),$$

where  $g(\omega) = 0$ , if  $\|\omega\|_2 \leq F$  and  $\infty$  otherwise. To see that, note that we can equivalently write the above as

$$\begin{aligned} \omega_{t+1} &= \arg \min_{\omega} \left\| \omega - \left( \omega_t - \frac{1}{L_2} \nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega_t}(\mathcal{D}_n) \right) \right\|_2, \text{ s.t. } \|\omega\|_2 \leq F \\ &= \text{proj}_{\omega: \|\omega\|_2 \leq F} \left( \omega_t - \frac{1}{L_2} \nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega_t}(\mathcal{D}_n) \right). \end{aligned}$$

Theorem 5 of (Karimi et al., 2016) gives us linear rates of convergence for projected gradient descent under the *proximal PL* condition. This condition is shown in Appendix G of (Karimi et al., 2016) to be equivalent to the following condition.

**Definition D.4.** A function  $F$  is said to satisfy the Kurdyka-Lojasiewicz condition with exponent  $1/2$  if there exists  $C > 0$  such that

$$\min_{s \in \partial F(\omega)} \|s\|_2^2 \geq C(F(\omega) - F_*),$$

where  $\partial F(\omega)$  is the Frechet subdifferential of  $F$  at  $\omega$  and  $F_*$  denotes the minimum value of  $F$ .

Note that, in our case we have  $F(\omega) = \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n) - g(\omega)$  and the Frechet subdifferential of this function in the domain  $\{\omega : \|\omega\|_2 \leq F\}$  only contains  $\nabla_{\omega} \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$ . Thus, the above condition is equivalent to the PL condition. As a consequence Theorem 5 of (Karimi et al., 2016) implies that

$$\mathcal{L}_{\text{RLHF}}^{\omega_t}(\mathcal{D}_n) - \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n) \leq \left(1 - \frac{C_{PL}}{L_2}\right)^t (\mathcal{L}_{\text{RLHF}}^{\omega_0}(\mathcal{D}_n) - \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n)).$$

Now, recalling the Hessian of our loss, note that for any non-zero vector  $v \in \mathbb{R}^{d_P}$ , we have

$$\begin{aligned} &v^{\top} \nabla_{\omega}^2 \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n) v \\ &= v^{\top} \left( \frac{1}{n} \sum_{(x, y^w, y^l) \sim \mathcal{D}_n} \frac{\exp(\omega^{\top} (\phi(x, y^w) - \phi(x, y^l)))}{(1 + \exp(\omega^{\top} (\phi(x, y^w) - \phi(x, y^l))))^2} (\phi(x, y^w) - \phi(x, y^l)) (\phi(x, y^w) - \phi(x, y^l))^{\top} \right) v \\ &\geq \frac{\exp(-F)}{(1 + \exp(2F))^2} \|v\|_{\Sigma_{\mathcal{D}_n, R}}^2. \end{aligned}$$

Thus,  $\mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n)$  is  $\frac{\exp(-F)}{(1 + \exp(2F))^2}$ -strongly convex with respect to the semi-norm  $\|\cdot\|_{\Sigma_{\mathcal{D}_n, R}}$  around  $\omega^*$ . Therefore, for any  $\omega$  we have

$$\begin{aligned} \mathcal{L}_{\text{RLHF}}^{\omega}(\mathcal{D}_n) - \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n) &\geq \langle \nabla_{\omega} \mathcal{L}_{\text{RLHF}}^*(\mathcal{D}_n), \omega - \omega^* \rangle + \frac{\exp(-2F)}{(1 + \exp(2F))^2} \|\omega - \omega^*\|_{\Sigma_{\mathcal{D}_n, R}}^2 \\ &\geq \frac{\exp(-2F)}{(1 + \exp(2F))^2} \|\omega - \omega^*\|_{\Sigma_{\mathcal{D}_n, R}}^2. \end{aligned}$$

Putting everything together, we obtain

$$\|\omega - \omega^*\|_{\Sigma_{\mathcal{D}_n, R}}^2 \leq O\left(\left(1 - \frac{1}{n}\right)^t\right).$$

□

## E. Convergence of Natural Policy Gradient for RLHF (Section 5)

In this section, we prove fast convergence rates for the a version of the natural policy gradient algorithm with loglinear policy class. We begin by deriving the gradient of the KL-regularized objective. Throughout, let us fix reward function  $r$  and dataset  $\mathcal{D}_n$ .

**Lemma E.1.** *Given reward  $r$  and policy  $\pi_\theta$ , the gradient of  $\mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n)$ , for the softmax policy class can be written as*

$$\nabla_\theta \mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n) = \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \bar{\psi}_\theta(x, y),$$

where

$$\bar{\psi}_\theta(x, y) = \psi(x, y) - \sum_{y' \in \mathcal{Y}} \pi_\theta(y'|x) \psi(x, y').$$

*Proof.* First, note that, for softmax policies with linear action preferences, we have, for any given  $(x, y)$ ,

$$\begin{aligned} \nabla_\theta \pi_\theta(y|x) &= \nabla_\theta \frac{\exp(\theta^\top \psi(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y'))} \\ &= \frac{\exp(\theta^\top \psi(x, y)) \sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y'))}{(\sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y')))^2} \psi(x, y) - \frac{\exp(\theta^\top \psi(x, y)) \sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y')) \psi(x, y')}{(\sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y')))^2} \psi(x, y) \\ &= \pi_\theta(y|x) (\psi(x, y) - \mathbb{E}_{y' \sim \pi_\theta(\cdot|x)}[\psi(x, y')]). \end{aligned}$$

On the other hand, for the regularizer, we have

$$\begin{aligned} \nabla_\theta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) &= \nabla_\theta (\theta^\top \psi(x, y) - \log \bar{Z}_\theta(x) - \log \mu(y|x)) \\ &= \psi(x, y) - \frac{1}{\bar{Z}_\theta(x)} \sum_{y' \in \mathcal{Y}} \exp(\theta^\top \psi(x, y')) \psi(x, y) \\ &= \psi(x, y) - \mathbb{E}_{y' \sim \pi_\theta(\cdot|x)}[\psi(x, y')], \end{aligned}$$

where

$$\bar{Z}_\theta(x) = \sum_{y \in \mathcal{Y}} \exp(\theta^\top \psi(x, y)).$$

Using the definition of  $\mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n)$  and the above derivations, we have

$$\begin{aligned} \nabla_\theta \mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n) &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \nabla_\theta \left( \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \right) \\ &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \left( \psi(x, y) - \sum_{y' \in \mathcal{Y}} \pi_\theta(y'|x) \psi(x, y') \right) \\ &\quad - \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( \psi(x, y) - \sum_{y' \in \mathcal{Y}} \pi_\theta(y', x) \psi(x, y') \right) \\ &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \left( \psi(x, y) - \sum_{y' \in \mathcal{Y}} \pi_\theta(y'|x) \psi(x, y') \right) \\ &\quad - \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \psi(x, y) + \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y' \in \mathcal{Y}} \pi_\theta(y'|x) \psi(x, y') \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \left( \psi(x, y) - \sum_{y' \in \mathcal{Y}} \pi_\theta(y'|x) \psi(x, y') \right) \\
 &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_\theta(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right) \bar{\psi}_\theta(x, y).
 \end{aligned}$$

□

Now we consider the following gradient update rule. Let  $\Psi_n \in \mathbb{R}^{d_F \times nY}$  denote the feature matrix corresponding to  $\mathcal{D}_n$  with columns  $\psi(x, y)$ , for every  $(x, y) \in \mathcal{D}_n \times \mathcal{Y}$ . We will assume that  $\Psi_n$  is full column rank. For every  $t \geq 0$ , let

$$\theta_{t+1} = \theta_t + \eta' (\Psi_n \Psi_n^\top)^\dagger \nabla_\theta \mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n). \quad (11)$$

First, we will expand this gradient expression in the following lemma.

**Lemma E.2.** *For every  $t \geq 0$ , the gradient update can be written as*

$$\theta_{t+1} = \theta_t + \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) \alpha_t,$$

where  $\mathbb{R}^{nY \times nY} \ni H(\pi) = \text{diag}(\pi) - M(\pi)$ , for any policy  $\pi$ , with  $M(\pi)$  being a block-diagonal matrix composed of  $n$  blocks  $\pi(x)\pi(x)^\top \in \mathbb{R}^{Y \times Y}$ , with  $\pi(x) = [\pi(y|x)]_{y \in \mathcal{Y}}$ , for every  $x \in \mathcal{D}_n$ , and

$$\alpha_t = \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1}.$$

Here,  $r = [r(x, y)]_{(x, y) \in \mathcal{D}_n \times \mathcal{Y}}$  denotes the reward vector, and  $\log \mu = [\log \mu(y|x)]_{(x, y) \in \mathcal{D}_n \times \mathcal{Y}}$  denotes the vector of log values for the reference policy.

*Proof.* Using the gradient update and Lemma E.1, we have

$$\begin{aligned}
 \theta_{t+1} &= \theta_t + \eta' (\Psi_n \Psi_n^\top)^\dagger \nabla_\theta \mathcal{V}_r^{\pi_\theta}(\mathcal{D}_n) \\
 &= \theta_t + \eta' (\Psi_n \Psi_n^\top)^\dagger \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \left( r(x, y) - \beta \log \left( \frac{\pi_{\theta_t}(y|x)}{\mu(y|x)} \right) \right) \bar{\psi}_{\theta_t}(x, y) \\
 &= \theta_t + \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \left( r(x, y) - \beta \theta_t^\top \psi(x, y) + \beta \log \bar{Z}_{\theta_t}(x) + \beta \log \mu(y|x) \right) \bar{\psi}_{\theta_t}(x, y) \\
 &= \theta_t + \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \left( r(x, y) - \beta \theta_t^\top \psi(x, y) + \beta \log \mu(y|x) \right) \bar{\psi}_{\theta_t}(x, y) \\
 &= \theta_t + \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) (r - \beta \Psi_n^\top \theta_t + \beta \log \mu) \\
 &= \theta_t - \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right),
 \end{aligned}$$

where the fourth equality follows from the observation that

$$\beta \sum_y \pi_\theta(y|x) \log \bar{Z}_\theta(x) \bar{\psi}_\theta(x, y) = \beta \log \bar{Z}_\theta(x) \sum_y \pi_\theta(x, y) (\psi(x, y) - \mathbb{E}_{y' \sim \pi_\theta(\cdot|x)}[\psi(x, y')]) = 0,$$

while the last equality follows from the fact that  $H(\pi_\theta) c \mathbf{1} = \mathbf{0}$ , for any constant  $c$ . □

Now we will express  $\alpha_{t+1}$  in a different way using the above derivation.

**Lemma E.3.** *For every  $t \geq 0$ , we have*

$$\alpha_{t+1} = (I - (\eta' \beta / n) H(\pi_{\theta_t})) \alpha_t$$



*Proof.* Observe that

$$\begin{aligned}
 \alpha_{t+1} &= \beta \Psi_n^\top \theta_{t+1} - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_{t+1} - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \\
 &= \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} + \beta \Psi_n^\top (\theta_{t+1} - \theta_t) - \frac{\beta (\Psi_n^\top (\theta_t - \theta_{t+1}))^\top \mathbf{1}}{Y} \cdot \mathbf{1} \\
 &= \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} - \frac{\beta (\Psi_n^\top (\theta_t - \theta_{t+1}))^\top \mathbf{1}}{Y} \cdot \mathbf{1} \\
 &\quad - \beta \Psi_n^\top \left( \frac{\eta'}{n} (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right) \right) \\
 &= \left( I - (\beta \eta' / n) \Psi_n^\top (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) \right) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right) \\
 &\quad - \frac{\beta (\Psi_n^\top (\theta_t - \theta_{t+1}))^\top \mathbf{1}}{Y} \cdot \mathbf{1} \\
 &= (I - (\beta \eta' / n) H(\pi_{\theta_t})) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right) \\
 &\quad - \frac{\beta (\Psi_n^\top (\theta_t - \theta_{t+1}))^\top \mathbf{1}}{Y} \cdot \mathbf{1},
 \end{aligned}$$

where the third equality uses Lemma E.2 and the last equality follows from the fact that

$$\Psi_n^\top (\Psi_n \Psi_n^\top)^\dagger \Psi_n = \Psi_n^\top (\Psi_n^\dagger)^\top \Psi_n^\dagger \Psi_n = (\Psi_n^\dagger \Psi_n)^\top \Psi_n^\dagger \Psi_n = I,$$

since we assume  $\Psi_n$  to be full column rank. For the last term in the derivation above, we have

$$\begin{aligned}
 &\frac{\beta (\Psi_n^\top (\theta_t - \theta_{t+1}))^\top \mathbf{1}}{Y} \cdot \mathbf{1} \\
 &= \frac{\beta \eta'}{nY} \left( \Psi_n^\top (\Psi_n \Psi_n^\top)^\dagger \Psi_n H(\pi_{\theta_t}) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right) \right)^\top \mathbf{1} \cdot \mathbf{1} \\
 &= \frac{\beta \eta'}{nY} \left( H(\pi_{\theta_t}) \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right) \right)^\top \mathbf{1} \cdot \mathbf{1} \\
 &= \frac{\beta \eta'}{nY} \left( \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \cdot \mathbf{1} \right)^\top H(\pi_{\theta_t})^\top \mathbf{1} \cdot \mathbf{1} \\
 &= \mathbf{0},
 \end{aligned}$$

where we have used the fact that  $H(\pi_\theta)^\top \mathbf{1} = \mathbf{0}$ . The result follows.  $\square$

Next, we will decompose the matrix  $H(\pi_\theta)$  into simpler pieces and explore its structure.

**Lemma E.4.** *The eigenvalues of  $H(\pi_\theta)$  satisfy the following. The lowest eigenvalue is  $\lambda_1 = 0$  with multiplicity  $n$  with corresponding eigenvectors  $e_i$ , for each  $i \in [n]$ , where  $e_i$  are the vectors of ones in indices  $Y(i-1)$  to  $Yi$  and zeros everywhere else. Furthermore, we have that  $\min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi(y|x) \leq \lambda_2$  and  $\lambda_{\max} \leq \max_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi(y|x)$ .*

*Proof.* Again, given  $x \in \mathcal{D}_n$ , let us denote by  $\pi(x)$  the vector  $[\pi(y|x)]_{y \in \mathcal{Y}}$ . Note that, for any policy  $\pi$ , we can write  $H(\pi)$  as

$$H(\pi) = \begin{bmatrix} \text{diag}(\pi(x_1)) & 0 & \dots & 0 \\ 0 & \text{diag}(\pi(x_2)) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \text{diag}(\pi(x_n)) \end{bmatrix}$$

$$\begin{aligned}
 & - \begin{bmatrix} \pi(x_1)\pi(x_1)^\top & 0 & \dots & 0 \\ 0 & \pi(x_2)\pi(x_2)^\top & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \pi(x_n)\pi(x_n)^\top \end{bmatrix} \\
 & = \begin{bmatrix} H(\pi(x_1)) & 0 & \dots & 0 \\ 0 & H(\pi(x_2)) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & H(\pi(x_n)) \end{bmatrix}
 \end{aligned}$$

Now, given  $x \in \mathcal{D}_n$ , Lemma 22 of (Mei et al., 2020) states that the spectrum of  $H(\pi(x))$  satisfies  $\lambda_1 = 0$  with corresponding eigenvector  $\mathbf{1} \in \mathbb{R}^Y$ , and

$$\pi(y_{i-1}|x) \leq \lambda_i \leq \pi(y_i|x),$$

for each  $2 \leq i \leq Y$ , where  $\lambda_1 \leq \dots \leq \lambda_Y$  and  $\pi(y_1|x) \leq \dots \leq \pi(y_Y|x)$ . Furthermore, it is known that the spectrum of a block diagonal matrix is composed of the eigenvalues of each block, counting multiplicities. Thus, we have that 0 is the lowest eigenvalue of  $H(\pi)$  occurring with multiplicity  $n$ . The rest follows.  $\square$

**Lemma E.5.** *Let  $v \in \mathbb{R}^{nY}$  be any given vector. Then, we have that*

$$\left\| (I - H(\pi)) \left( v - \frac{v^\top \mathbf{1}}{Y} \mathbf{1} \right) \right\|_2 \leq \left( 1 - \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi(y|x) \right) \left\| v - \frac{v^\top \mathbf{1}}{Y} \mathbf{1} \right\|_2$$

*Proof.* First, for every  $i \in [n]$ , let  $e_i \in \mathbb{R}^{nY}$  denote a vector with entries 1 at the indices  $Y(i-1)$  to  $Yi$  and 0 everywhere else. Note that

$$\sum_{i \leq n} e_i = \mathbf{1} \in \mathbb{R}^{nY}.$$

Next, let  $v(j)$  denote an  $nY$ -dimensional vector with entries  $v_k$  for each  $Y(j-1) \leq k \leq Yj$ . Since  $H(\pi)$  is diagonalizable, as a symmetric matrix, any vector can be represented as a linear combination of its eigenvectors. Since  $H(\pi)$  is symmetric, this representation is unique. Now, by Lemma E.4, note that

$$v = \sum_{j \leq nY} a_j u_j = \sum_{j \leq n} a_j e_j + \sum_{k=n+1}^{nY} a_k u_k = \sum_{j \leq n} \frac{v(j)^\top \mathbf{1}_Y}{Y} e_j + \sum_{k=n+1}^{nY} a_k u_k = \frac{v^\top \mathbf{1}}{Y} \mathbf{1} + \sum_{k=n+1}^{nY} a_k u_k,$$

where  $(u_{i,j})_{i \leq n, j \leq Y}$  is the eigenvector basis, with the first  $n$  eigenvectors being  $e_i$ , for  $i \leq n$ . Thus, we have that

$$v' = v - \frac{v^\top \mathbf{1}}{Y} \mathbf{1} = \sum_{k=n+1}^{nY} a_k u_k,$$

with  $a_{n+1} > 0$ , and that

$$\|v'\|_2 = \sum_{j=n+1}^{nY} a_j^2.$$

From the above, we obtain

$$(I - H(\pi)) v' = \sum_{j=n+1}^{nY} a_j (1 - \lambda_j) u_j,$$

and thus, by Lemma E.4,

$$\|(I - H(\pi)) v'\|_2 = \sqrt{\sum_{j=n+1}^{nY} a_j^2 (1 - \lambda_j)^2}$$

$$\begin{aligned}
 &\leq \sqrt{(1 - \lambda_{n+1}) \left( \sum_{j=n+1}^{nY} a_j^2 \right)} \\
 &= (1 - \lambda_{n+1}) \|v'\|_2 \\
 &\leq \left( 1 - \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi(y|x) \right) \|v'\|_2.
 \end{aligned}$$

□

**Lemma E.6.** Suppose  $\eta'$  and  $\beta$  are such that  $\eta' \beta / n \leq 1$ . For the loglinear policy class, for every  $t \geq 1$ ,

$$\|\alpha_t\|_2 \leq \frac{2(\beta B + 1) \sqrt{Y}}{\exp\left(\eta' \beta \sum_{s=1}^{t-1} \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_s}(y|x)\right)}.$$

*Proof.* By Lemma E.3 and Lemma E.5, for all  $t \geq 1$ ,

$$\begin{aligned}
 \|(I - (\eta' \beta / n) H(\pi_{\theta_{t+1}})) \alpha_{t+1}\|_2 &\leq \left( 1 - (\eta' \beta / n) \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \right) \|\alpha_t\|_2 \\
 &\leq \frac{1}{\exp\left((\eta' \beta / n) \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x)\right)} \|\alpha_t\|_2 \\
 &\leq \frac{1}{\exp\left((\eta' \beta / n) \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x)\right)} \left( 1 - \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_{t-1}}(y|x) \right) \|\alpha_{t-1}\|_2 \\
 &\leq \frac{1}{\exp\left((\eta' \beta / n) \sum_{s=t-1}^t \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_s}(y|x)\right)} \|\alpha_{t-1}\|_2 \\
 &\leq \frac{1}{\exp\left((\eta' \beta / n) \sum_{s=1}^t \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_s}(y|x)\right)} \|\alpha_1\|_2.
 \end{aligned}$$

For the first iteration, observe that

$$\begin{aligned}
 \|\alpha_1\|_2 &= \left\| \beta \Psi_n^\top \theta_0 - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_0 - r - \beta \log \mu)^\top \mathbf{1}}{Y} \mathbf{1} \right\|_2 \\
 &\leq \|\beta \Psi_n^\top \theta_0 - r - \beta \log \mu\|_2 + \frac{1}{\sqrt{Y}} \|\beta \Psi_n^\top \theta_0 - r - \beta \log \mu\|_2 \|\mathbf{1}\|_2 \\
 &= 2 \|\beta \Psi_n^\top \theta_0 - r - \beta \log \mu\|_2 \\
 &\leq 2(\beta \|\Psi_n^\top \theta_0\|_\infty + 1) \sqrt{Y} \\
 &\leq 2(\beta B + 1) \sqrt{Y},
 \end{aligned}$$

where the second inequality follows from the triangle inequality, the Cauchy-Schwarz inequality and the fact that rewards lie in the unit ball, while the last follows from the fact that the features lie in a unit subspace of  $\mathbb{R}^{d_p}$ , while  $\|\theta_0\|_\infty \leq B$ . The result follows. □

**Lemma E.7.** There exists a constant  $C = C(\beta, Y, B) > 0$ , such that, for all  $t \geq 1$ , we have  $\min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \geq C$ .

*Proof.* First, by Lemma E.5, note that, for any  $t \geq 1$ ,

$$\|\alpha_{t+1}\|_2 \leq \left( 1 - (\eta' \beta / n) \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x) \right) \|\alpha_t\|_2 \leq \|\alpha_t\|_2 \leq \dots \leq \|\alpha_1\|_2 \leq 2(\beta B + 1) \sqrt{Y},$$

where the second inequality follows from the fact that policies are probability distributions, and the last inequality follows from Lemma E.6. Next, observe that, for any  $(x, y) \in \mathcal{D}_n \times \mathcal{Y}$ ,

$$\left| \psi(x, y)^\top \theta_t - \frac{1}{\beta} r(x, y) - \log \mu(y|x) - \frac{(\Psi_n^\top \theta_t - r / \beta - \log \mu)^\top \mathbf{1}}{Y} \right|$$

$$\begin{aligned}
 &\leq \frac{1}{\beta} \left| \beta \psi(x, y)^\top \theta_t - r(x, y) - \beta \log \mu(y|x) - \frac{\beta (\Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \right| \\
 &\leq \frac{1}{\beta} \left\| \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \mathbf{1} \right\|_2 \\
 &\leq \frac{1}{\beta} \|\alpha_t\|_2 \\
 &\leq 2(B + 1/\beta) \sqrt{Y}.
 \end{aligned}$$

Now, define  $(x_1, y_1) = \arg \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \psi(x, y)^\top \theta_t$  and  $(x_2, y_2) = \arg \max_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \psi(x, y)^\top \theta_t$ . By the above, we have

$$\begin{aligned}
 \Psi_n(x_1, y_1)^\top \theta_t &\geq \frac{1}{\beta} r(x_1, y_1) + \log \mu(y_1|x_1) + \frac{(\Psi_n^\top \theta_t - r/\beta - \log \mu)^\top \mathbf{1}}{Y} - 2(B + 1/\beta) \sqrt{Y}, \\
 -\Psi_n(x_2, y_2)^\top \theta_t &\geq -\frac{1}{\beta} r(x_2, y_2) - \log \mu(y_2|x_2) - \frac{(\Psi_n^\top \theta_t - r/\beta - \log \mu)^\top \mathbf{1}}{Y} - 2(B + 1/\beta) \sqrt{Y},
 \end{aligned}$$

which imply

$$\begin{aligned}
 \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \pi_{\theta_t}(y|x) &\geq \min_{x \in \mathcal{D}_n, y \in \mathcal{Y}} \frac{\exp(\psi(x, y)^\top \theta_t)}{\sum_{y' \in \mathcal{Y}} \exp(\psi(x, y')^\top \theta_t)} \geq \frac{1}{Y} \exp\left((\Psi_n(x_1, y_1) - \Psi_n(x_2, y_2))^\top \theta_t\right) \\
 &\geq \frac{1}{Y} \exp\left(\frac{1}{\beta} (r(x_1, y_1) - r(x_2, y_2)) + \log \frac{\mu(y_1|x_1)}{\mu(y_2|x_2)} - 4(B + 1/\beta) \sqrt{Y}\right) \\
 &\geq \frac{1}{Y} \exp\left(-\frac{1}{\beta} - 4(B + 1/\beta) \sqrt{Y}\right) = C.
 \end{aligned}$$

□

Let us denote by  $\text{softmax}(\Psi_n^\top v)$  the policy  $\exp(\psi(x, y)^\top v) / \sum_{y'} \exp(\psi(x, y')^\top v)$ , for any parameter vector  $v$  and pair  $(x, y) \in \mathcal{D}_n \times \mathcal{Y}$ . Now, we are ready to prove the main result of this section.

**Theorem E.8.** *Let  $\pi_{\theta_t} = \text{softmax}(\Psi_n^\top \theta_t)$ . Using update rule (11) with  $\eta' \leq n/\beta$ , for all  $t \geq 1$ ,*

$$\mathcal{V}_r^{\pi_{\theta^*}}(\mathcal{D}_n) - \mathcal{V}_r^{\pi_{\theta_t}}(\mathcal{D}_n) \leq \frac{2\sqrt{Y}(B + 1/\beta)}{\exp((\beta\eta'/n) \cdot C \cdot (t - 1))},$$

where

$$C = \frac{1}{Y} \exp\left(-\frac{1}{\beta} - 4(B + 1/\beta) \sqrt{Y}\right).$$

*Proof.* Observe that, since  $\pi_{\theta^*} \propto \mu(y|x) \exp(r(x, y)/\beta)$ , we have

$$\begin{aligned}
 \mathcal{V}_r^{\pi_{\theta^*}}(\mathcal{D}_n) - \mathcal{V}_r^{\pi_{\theta_t}}(\mathcal{D}_n) &= \frac{1}{n} \sum_{x \in \mathcal{D}_n} \sum_{y \in \mathcal{Y}} (\pi_{\theta^*}(y|x) r(x, y) - \beta D_{\text{KL}}(\pi_{\theta^*} \|\mu) - \pi_{\theta_t}(y|x) r(x, y) + \beta D_{\text{KL}}(\pi_{\theta_t} \|\mu)) \\
 &\leq \frac{1}{n} \sum_{x \in \mathcal{D}_n} \|\pi_{\theta^*}(\cdot|x) - \pi_{\theta_t}(\cdot|x)\|_1 \\
 &\quad + \frac{1}{n} \sum_{x \in \mathcal{D}_n} -\beta D_{\text{KL}}(\pi_{\theta^*} \|\pi_{\theta^*}) + \pi_{\theta^*}(y|x) r(x, y) + \beta D_{\text{KL}}(\pi_{\theta_t} \|\pi_{\theta^*}) - \pi_{\theta_t}(y|x) r(x, y) \\
 &\leq \frac{2}{n} \sum_{x \in \mathcal{D}_n} \|\pi_{\theta^*}(\cdot|x) - \pi_{\theta_t}(\cdot|x)\|_1 + \beta D_{\text{KL}}(\pi_{\theta_t} \|\pi_{\theta^*}) \\
 &\leq (2Y + \beta) D_{\text{KL}}(\pi_{\theta_t} \|\pi_{\theta^*})
 \end{aligned}$$

$$\leq (2Y + \beta) \left\| \Psi_n^\top \theta^* - \Psi_n \theta_t + \frac{(\Psi_n^\top (\theta_t - \theta^*))^\top \mathbf{1}}{Y} \mathbf{1} \right\|_\infty^2,$$

where the third inequality uses Pinsker's inequality and the last one follows from Lemma J.8. Now, note that the optimal softmax policy parameter  $\theta^*$  satisfies, for each  $(x, y) \in \mathcal{D}_n$ ,

$$\psi(x, y)^\top \theta^* = \frac{1}{\beta} (r(x, y) + \log \mu(y|x)),$$

by setting the gradient at  $(x, y)$  to 0. Its existence is guaranteed by the assumption that  $r^* \in \mathcal{F}$  and Lemma J.1. Thus, we have

$$\begin{aligned} \mathcal{V}_r^{\pi_{\theta^*}}(\mathcal{D}_n) - \mathcal{V}_r^{\pi_{\theta_t}}(\mathcal{D}_n) &\leq (2Y + \beta) \left\| \Psi_n^\top \theta^* - \Psi_n \theta_t + \frac{(\Psi_n^\top (\theta_t - \theta^*))^\top \mathbf{1}}{Y} \mathbf{1} \right\|_\infty^2 \\ &= (2Y + \beta) \left\| \frac{1}{\beta} (r + \log \mu) - \Psi_n^\top \theta_t + \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{\beta Y} \mathbf{1} \right\|_\infty^2 \\ &= \frac{(2Y + \beta)}{\beta} \left\| \beta \Psi_n^\top \theta_t - r - \beta \log \mu - \frac{(\beta \Psi_n^\top \theta_t - r - \beta \log \mu)^\top \mathbf{1}}{Y} \mathbf{1} \right\|_\infty^2 \\ &\leq \frac{2(2Y + \beta)2\sqrt{Y}(\beta B + 1)}{\beta \exp((\eta' \beta/n)C(t-1))}, \end{aligned}$$

where the last inequality follows from Lemma E.6 and Lemma E.7.  $\square$

## F. Convergence of Gradient Descent for DPO (Section 5)

In this section, we will prove convergence bounds for the projected gradient descent procedure for DPO. Recall that the projected gradient descent is defined as

$$\theta_{t+1} = \underset{\theta: \|\theta\|_2 \leq B}{\text{proj}} (\theta_t - \eta' \nabla_\theta \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)),$$

We begin by showing that the DPO objective satisfies the PL condition (Karimi et al., 2016) stated in Definition D.2. We will show that the DPO objective satisfies this condition for the loglinear parametrization. First, we need to show that such an objective has Lipschitz gradients, which holds under the assumption that the parameter vectors  $\theta$  have a length of no more than  $B$ .

**Lemma F.1.** *The DPO objective  $\mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)$  is Lipschitz continuous with parameter  $L'_1 = \beta \exp(2\beta(B + J))$  and has Lipschitz gradients with parameter  $L'_2 = \beta^2 \exp(2\beta(B + J))$ , where*

$$J = \max_{(x, y^w, y^l) \in \mathcal{D}_n} \beta \left| \log \frac{\mu(y^w|x)}{\mu(y^l|x)} \right|.$$

*Proof.* Note that, in order to show  $L$ -Lipschitzness, it suffices to prove that the Hessian of  $\mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)$  has bounded eigenvalues. Let us first compute the Hessian. Before doing that, we first simplify the gradient expression, when instantiated for the softmax parametrization. First, given parameter vector  $\theta$ , corresponding to  $\pi_\theta$ , we have

$$\begin{aligned} \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n) &= -\mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}_n} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y^w|x)}{\mu(y^w|x)} - \beta \log \frac{\pi_\theta(y^l|x)}{\mu(y^l|x)} \right) \right] \\ &= -\mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}_n} \left[ \log \sigma \left( \beta \log \frac{\exp(\theta^\top \psi(x, y^w))}{\sum_{y \in \mathcal{Y}} \exp(\theta^\top \psi(x, y))} - \beta \log \frac{\exp(\theta^\top \psi(x, y^l))}{\sum_{y \in \mathcal{Y}} \exp(\theta^\top \psi(x, y))} - \beta \log \frac{\mu(y^w|x)}{\mu(y^l|x)} \right) \right] \\ &= -\mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}_n} \left[ \log \sigma \left( \beta \theta^\top (\psi(x, y^w) - \psi(x, y^l)) - \beta \log \frac{\mu(y^w|x)}{\mu(y^l|x)} \right) \right] \\ &= \mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}_n} \left[ \log (1 + \exp(\beta \theta^\top (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l))) \right], \end{aligned}$$

where we let

$$J(x, y^w, y^l) = \beta \log \frac{\mu(y^w|x)}{\mu(y^l|x)}.$$

Based on the above, we have

$$\begin{aligned} \nabla_{\theta} \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n) &= \nabla_{\theta} \mathbb{E}_{(y^w, y^l, x) \sim \mathcal{D}_n} [\log(1 + \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l)))] \\ &= \frac{1}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \frac{\beta \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l))}{(1 + \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l)))} (\psi(x, y^w) - \psi(x, y^l)), \end{aligned}$$

and

$$\begin{aligned} \nabla_{\theta}^2 \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n) &= \frac{1}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \nabla_{\theta} \frac{\beta \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l))}{(1 + \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l)))} (\psi(x, y^w) - \psi(x, y^l)) \\ &= \frac{1}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \frac{\beta^2 \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l))}{(1 + \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l)))^2} (\psi(x, y^w) - \psi(x, y^l)) (\psi(x, y^w) - \psi(x, y^l))^{\top}. \end{aligned}$$

Now, define

$$E(\theta, x, y) = \exp(\beta \theta^{\top} (\psi(x, y^w) - \psi(x, y^l)) - J(x, y^w, y^l)).$$

Note that we have

$$\|\nabla \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)\|_2 \leq \beta \exp(2\beta(B + J)) \|\psi(x, y^w) - \psi(x, y^l)\|_2 \leq \beta \exp(2\beta(B + J)),$$

and

$$\begin{aligned} \nabla_{\theta}^2 \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n) &= \beta^2 \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \frac{E(\theta, x, y)}{n(1 + E(\theta, x, y))^2} \psi(x) \psi(x)^{\top} \\ &\preceq \frac{\beta^2 \exp(2\beta(B + J))}{n} \sum_{(x, y^w, y^l) \in \mathcal{D}_n} \psi(x) \psi(x)^{\top} \\ &\preceq \beta^2 \exp(2\beta(B + J)) I_d, \end{aligned}$$

where the last inequality follows from the fact that the feature norms are bounded by 1, and thus the maximum eigenvalue of the sample covariance matrix is no more than 1.  $\square$

Next, we show that the DPO objective satisfies the PL condition under some mild assumption on the data.

**Lemma F.2.** *Assume that, for each triple  $(x, y^w, y^l) \in \mathcal{D}_n$ , we have that  $\psi(x, y^w) \neq \psi(x, y^l)$ . Then, if we let*

$$C'_{PL} = \frac{\beta \exp(-2\beta(B + J))^3 (1 + \exp(-2\beta(B + J)))}{n(1 + \exp(2\beta(B + J)))^2} \min_{(x, y^w, y^l) \in \mathcal{D}_n} \|\psi(x, y^w) - \psi(x, y^l)\|^2,$$

we have

$$\frac{1}{2} \|\nabla \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)\|^2 \geq C'_{PL} (\mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n) - \mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n))$$

where  $\mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n) = \min_{\theta} \mathcal{L}_{\text{DPO}}^{\theta}(\mathcal{D}_n)$  denotes the optimal loss value.

*Proof.* Using the notation  $E(\theta, x, y) = \exp(\beta\theta^\top(\psi(x, y^w) - \psi(x, y^l)) - E(x))$ , and noting that every quantity in the expression below is non-negative, we have

$$\begin{aligned} \frac{1}{2} \|\nabla_\theta \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)\|^2 &\geq \frac{\beta^2}{2n^2} \sum_{(x,y) \in \mathcal{D}_n} \frac{E(\theta, x, y)^2}{(1 + E(\theta, x, y))^2} \|\psi(x, y^w) - \psi(x, y^l)\|^2 \\ &\geq \frac{\beta^2 \exp(-2\beta(B + J))^2}{n(1 + \exp(2\beta(B + J)))^2} \min_{(x,y) \in \mathcal{D}_n} \|\psi(x, y^w) - \psi(x, y^l)\|^2 \end{aligned}$$

since  $-J \leq E(x) \leq J$ , and thus,  $\exp(-2\beta(B + J)) \leq E(\theta, x, y) \leq \exp(2\beta(B + J))$ . On the other hand, observe that

$$\begin{aligned} \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n) - \mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n) &= \mathbb{E}_{(x,y) \sim \mathcal{D}_n} \left[ \log(1 + E(\theta, x, y)) - \log(1 + E(\theta^*, x, y)) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \log \frac{1 + E(\theta, x_i, y_i)}{1 + E(\theta^*, x_i, y_i)} \\ &\leq \frac{1}{n} \sum_{i=1}^n \left( \frac{1 + E(\theta, x_i, y_i)}{1 + E(\theta^*, x_i, y_i)} - 1 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{E(\theta, x_i, y_i) - E(\theta^*, x_i, y_i)}{1 + E(\theta^*, x_i, y_i)} \\ &\leq \frac{1}{n} \sum_{i=1}^n \frac{E(\theta, x_i, y_i)}{1 + E(\theta^*, x_i, y_i)} \\ &\leq \frac{\exp(2\beta(B + J))}{1 + \exp(-2\beta(B + J))}. \end{aligned}$$

Now, the assumption on the data implies that, there exists  $\xi$  such that

$$0 < \xi' = \min_{(x,y^w,y^l) \in \mathcal{D}_n} \|\psi(x, y^w) - \psi(x, y^l)\|^2.$$

Using  $\xi'$  and solving the equation

$$C'_{PL} \cdot \frac{\exp(2\beta(B + J))}{1 + \exp(-2\beta(B + J))} = \frac{\beta \exp(-2\beta(B + J))^2}{n(1 + \exp(2\beta(B + J)))^2} \xi'$$

for  $C'_{PL}$ , we obtain

$$C'_{PL} = \frac{\beta \exp(-2\beta(B + J))^3 (1 + \exp(-2\beta(B + J)))}{n(1 + \exp(2\beta(B + J)))^2} \xi'. \quad (12)$$

□

These two conditions are enough to obtain the following result.

**Theorem 5.3.** *For every  $t \geq 0$ , the gradient descent procedure (6) with learning rate  $\eta'' = O(1/\beta^2)$  satisfies*

$$\|\theta_t - \theta_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, P}}^2 \leq O\left(\frac{1}{\beta} \left(1 - \frac{\beta}{n}\right)^t\right).$$

*Proof.* A similar argument as the one in the proof of Theorem 5.1 implies that, for every  $t \geq 1$  we have:

$$\mathcal{L}_{\text{DPO}}^{\theta_t}(\mathcal{D}_n) - \mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n) \leq \left(1 - \frac{C'_{PL}}{L'_2}\right)^t \left(\mathcal{L}_{\text{DPO}}^{\theta_0}(\mathcal{D}_n) - \mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n)\right),$$

where  $L'_1 = \beta \exp(2\beta(B + J))$  is the Lipschitz constant and  $C'_{LPL} = BC_{PL}$ . Using the expression for the Hessian derived in the proof of Lemma F.1 we have, for any non-zero vector  $v$ , that

$$v^\top \nabla_\theta^2 \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n)v \geq \beta^2 \frac{\exp(-\beta(B + J))}{1 + \exp(\beta(B + J))} \|v\|_{\Sigma_{\mathcal{D}_n, P}}^2.$$

Thus,  $\mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n)$  is  $\beta^2 \frac{\exp(-\beta(B+J))}{1+\exp(\beta(B+J))}$ -strongly convex with respect to the semi-norm  $\|\cdot\|_{\Sigma_{\mathcal{D}}}$  around  $\theta_{\mathcal{D}_n}^*$ , where  $\theta_{\mathcal{D}_n}^*$  is a parameter vector that achieves  $\mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n)$ . for any  $\theta$ , we have

$$\begin{aligned} \mathcal{L}_{\text{DPO}}^\theta(\mathcal{D}_n) - \mathcal{L}_{\mathcal{D}_n}^* &\geq \langle \nabla_\theta \mathcal{L}_{\text{DPO}}^*(\mathcal{D}_n), \theta - \theta_{\mathcal{D}_n}^* \rangle + \beta^2 \frac{\exp(-\beta(B + J))}{2(1 + \exp(\beta(B + J)))} \|\theta - \theta_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, P}}^2 \\ &\geq \beta^2 \frac{\exp(-\beta(B + J))}{2(1 + \exp(\beta(B + J)))} \|\theta - \theta_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, P}}^2 \end{aligned}$$

Therefore, using the upper bound on the loss, we finally obtain, for any iterate  $\theta_t$  of GD,

$$\begin{aligned} \|\theta_t - \theta_{\mathcal{D}_n}^*\|_{\Sigma_{\mathcal{D}_n, P}}^2 &\leq O\left(\frac{\mathcal{L}_{\text{DPO}}^{\theta_0}(\mathcal{D}_n) - \mathcal{L}_{\mathcal{D}_n}^*}{\beta^2} \left(1 - \frac{\beta}{n}\right)^t\right) \\ &\leq O\left(\frac{1}{\beta} \left(1 - \frac{\beta}{n}\right)^t\right) \end{aligned}$$

□

## G. Non-realizable Rewards (Section 6)

In this section, we will derive the proofs of the two results from Section 6. We restate them for convenience.

**Theorem 6.1.** *Let  $\delta > 0$ . Suppose that Assumption 6.1 holds. Then, with probability at least  $1 - \delta$ , we have  $G(\pi_{\hat{\theta}}) \leq D(\pi_{\hat{\theta}}) + \tilde{\Theta}(\Lambda_R \sqrt{d_R/n}) + 2\epsilon_{\text{app}}$ .*

*Proof.* From Theorem 4.1, we have

$$\begin{aligned} G(\pi_{\hat{\theta}}) &= D(\pi_{\hat{\theta}}) + \langle d_\rho^* - d_\rho^{\pi_{\hat{\theta}}}, r^* - r_{\hat{\omega}} \rangle \\ &= D(\pi_{\hat{\theta}}) + \langle d_\rho^* - d_\rho^{\pi_{\hat{\theta}}}, r^* - r_{\omega^*} \rangle + \langle d_\rho^* - d_\rho^{\pi_{\hat{\theta}}}, r_{\omega^*} - r_{\hat{\omega}} \rangle \\ &\leq D(\pi_{\hat{\theta}}) + 2 \max_{x,y} |r^*(x,y) - r_{\omega^*}(x,y)| + O\left(\Lambda_R \sqrt{\frac{d_R}{n}}\right) \\ &\leq D(\pi_{\hat{\theta}}) + O\left(\Lambda_R \sqrt{\frac{d_R}{n}}\right) + 2\epsilon_{\text{app}}, \end{aligned}$$

where for the first inequality we have used Cauchy-Schwarz, while for the last inequality we have used Theorem 4.1 and Condition 6.1. □

Next, we prove the analogous result for DPO.

**Theorem 6.2.** *Let  $\delta > 0$ . Suppose that Assumption 6.1 and the condition of Lemma 4.1 hold. Then, with probability at least  $1 - \delta$ , we have  $G(\pi_{\hat{\theta}}) \leq D(\pi_{\hat{\theta}}) + \Theta(\Lambda_P d_P / (\beta n)) + \min\{2\epsilon_{\text{app}}, O(\beta D_{\text{KL}}(\pi_{\theta^*} \|\pi^*))\}$ .*

*Proof.* Since the ground-truth reward function is not linear, we are not guaranteed that the optimal policy representable in terms of the reward is loglinear. Let  $\pi^*$  denote the optimal policy for the KL-regularized problem with respect to  $r^*$ , and let  $\pi_{\theta^*}$  be the loglinear approximation of  $\pi^*$ .

$$G(\pi_{\hat{\theta}}) = V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_{\hat{\theta}}}(\rho)$$



$$\begin{aligned}
 &= D(\pi_{\tilde{\theta}}) + \left( \mathcal{V}_{r^*}^{\pi^*}(\rho) - \mathcal{V}_{r^*}^{\pi_{\tilde{\theta}}}(\rho) \right) \\
 &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{x \sim \rho, y \sim \pi^*(\cdot|x)} \left[ r^*(x, y) - \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} \right] - \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ r^*(x, y) - \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} \right] \\
 &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{x \sim \rho, y \sim \pi^*(\cdot|x)} \left[ \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} + \beta \log Z(x) - \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} \right] \\
 &\quad - \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} + \beta \log Z(x) - \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} \right] \\
 &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} - \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} \right] \\
 &= D(\pi_{\tilde{\theta}}) + \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \beta \log \frac{\pi_{\tilde{\theta}}(y|x)}{\mu(y|x)} - \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} \right] \\
 &\quad + \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} \left[ \beta \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} - \beta \log \frac{\pi^*(y|x)}{\mu(y|x)} \right] \\
 &= D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + \mathbb{E}_{x \sim \rho, y \sim \pi_{\tilde{\theta}}(\cdot|x)} [\beta \log \pi_{\theta^*}(y|x) - \beta \log \pi^*(y|x)] \\
 &= D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + \beta \sum_x \rho(x) \sum_y \frac{\pi_{\tilde{\theta}}(y|x)}{\pi_{\theta^*}(y|x)} \log \frac{\pi_{\theta^*}(y|x)}{\pi^*(y|x)} \\
 &\leq D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + \beta Y \exp(2B) D_{\text{KL}}(\pi_{\theta^*} || \pi^*) .
 \end{aligned}$$

On the other hand, using the same idea as in the proof of Theorem 4.2, we have

$$\begin{aligned}
 G(\pi_{\tilde{\theta}}) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_{\tilde{\theta}}}(\rho) \\
 &= D(\pi_{\tilde{\theta}}) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\tilde{\theta}}}(\rho)) \\
 &= D(\pi_{\tilde{\theta}}) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r_{\omega^*}}^{\pi_{\theta^*}}(\rho)) + (\mathcal{V}_{r_{\omega^*}}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r_{\omega^*}}^{\pi_{\tilde{\theta}}}(\rho)) + (\mathcal{V}_{r_{\omega^*}}^{\pi_{\tilde{\theta}}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\tilde{\theta}}}(\rho)) \\
 &= D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} [r^*(x, y) - r_{\omega^*}(x, y)] + \mathbb{E}_{(x,y) \sim d_{\rho}^{\tilde{\theta}}} [r_{\omega^*}(x, y) - r^*(x, y)] \\
 &\leq D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + \mathbb{E}_{(x,y) \sim d_{\rho}^{\theta^*}} [r^*(x, y) - r_{\omega^*}(x, y)] + 2 \max_{x,y} |r^*(x, y) - r_{\omega^*}(x, y)| \\
 &\leq D(\pi_{\tilde{\theta}}) + \Theta \left( \frac{\Lambda(d_P + 1)}{\beta n} \right) + 2\epsilon_{\text{app}} ,
 \end{aligned}$$

where the fourth equality follows from Theorem 4.2 and the last inequality from Condition 6.1.  $\square$

## H. The DPO Extension to MDPs (Section 7)

First, we start with MDP setting preliminaries.

### H.1. Deterministic Markov Decision Processes

An infinite-horizon discounted deterministic Markov decision process (MDP) is a mathematical object  $\mathcal{M} = (\mathcal{X}, \mathcal{Y}, T, r^*, \gamma, \rho)$ , where  $\mathcal{X}$  denotes the state space,  $\mathcal{Y}$  denotes the action space, both of which are assumed to be finite with cardinalities  $X$  and  $Y$ , respectively.  $T : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{X}$  denotes the deterministic transition function, where  $T(x, y)$  denotes the next state after taking action  $y$  in state  $x$ . The reward function is denoted by  $r^* : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ . Finally,  $\gamma \in [0, 1)$  denotes the discount factor, while  $\rho \in \Delta(\mathcal{X})$  denotes the initial state distribution.

Policies  $\pi$  are mappings from states to distributions over actions, that is,  $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ . Given policy  $\pi$ , the state occupancy measure of state  $x$  with respect to initial state  $x_0$  is given as

$$d_{x_0}^{\pi}(x) = (1 - \gamma) \sum_{t \geq 0} \gamma^t \mathbb{P}(x_t = x | x_0, \pi) ,$$

while the state-action occupancy measure is given as  $d_{x_0}^\pi(x, y) = d_{x_0}^\pi(x)\pi(y|x)$ . We also write  $d_\rho^\pi(x, y) = \mathbb{E}_{x_0 \sim \rho}[d_{x_0}^\pi(x, y)]$ . Furthermore, given policy  $\pi$  and an arbitrary reward function  $r$ , the value function of policy  $\pi$  with respect to reward  $r$  is defined as

$$V_r^\pi(x) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(x_t, y_t) \middle| x_0 = x, \pi \right],$$

and the action-value function is defined as

$$Q_r^\pi(x, y) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r(x_t, y_t) \middle| x_0 = x, y_0 = y, \pi \right],$$

for every state-action pair  $(x, y)$ . We denote by  $V_r^\pi(\rho) = \mathbb{E}_{x \sim \rho}[V_r^\pi(x)]$  the expected value function over the initial distribution.

## H.2. DPO for MDPs

A direct extension of DPO to the MDP setting is not straightforward. To understand this, it is enough to see that the optimal policy-to-reward mapping in this case is not linear. Fix a reward function  $r$ . The gradient of the KL-regularized objective with respect to  $r$  is given as

$$\nabla_\theta \mathcal{V}_r^\theta(\rho) = \frac{1}{1-\gamma} \sum_x d_\rho^{\pi_\theta}(x) \sum_y \pi_\theta(y|x) \left( r(x, y) + \gamma \mathcal{V}_r^{\pi_\theta}(T(x, y)) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \bar{\psi}_\theta(x, y),$$

where  $\bar{\psi}_\theta(x, y) = \psi(x, y) - \mathbb{E}_{y' \sim \pi_\theta(\cdot|x)}[\psi(x, y')]$ . See Appendix I for derivations. What complicates things is the occupancy measure  $d_\rho^\pi$ , which is non-linearly dependent on policy  $\pi$ , and the gradient of the occupancy measure. To allow for the change of variables to carry through in this case, we utilize the dual formulation of Problem (P3.2):

$$\begin{aligned} \max_{d_\rho} \quad & \sum_{x,y} d_\rho(x, y) r(x, y) - \beta \sum_{x,y} d_\rho(x, y) \log \frac{d_\rho(x, y)}{d_\rho^\mu(x, y)} \\ \text{s.t.} \quad & \sum_y d_\rho(x, y) = (1-\gamma)\rho(x) + \gamma \sum_{x',y'} \mathbb{1}(x = T(x', y')) d_\rho(x', y'), \forall x \in \mathcal{X}, \end{aligned} \quad (\text{P3.2}')$$

where we have used that

$$V_r^\pi(\rho) = \sum_{x,y} d_\rho^\pi(x, y) r(x, y),$$

and also taken the KL-divergence of the occupancy measures, instead of the actual policies. This is a convex program and thus any stationary points are optimal. The Lagrangian of the above problem can be written as

$$\begin{aligned} L(d_\rho, \alpha) &= \sum_{x,y} d_\rho(x, y) \left( r^*(x, y) - \beta \log \frac{d_\rho(x, y)}{d_\rho^\mu(x, y)} \right) \\ &\quad + \sum_x \alpha(x) \left( \sum_y d_\rho(x, y) - (1-\gamma)\rho(x) - \gamma \sum_{x',y'} \mathbb{1}(x = T(x', y')) d_\rho(x', y') \right) \\ &= -\beta \sum_{x,y} d_\rho(x, y) \log \frac{d_\rho(x, y)}{d_\rho^\mu(x, y)} - (1-\gamma) \sum_x \rho(x) \alpha(x) \\ &\quad + \sum_{x,y} d_\rho(x, y) \underbrace{\left( r^*(x, y) - \gamma \sum_{x'} \mathbb{1}(x = T(x', y')) \alpha(x') + \alpha(x) \right)}_{e_\alpha(x,y)}, \end{aligned}$$

Then, given  $(x, y)$ , the gradient of the Lagrangian with respect to  $d_\rho(x, y)$  is

$$\nabla_{d_\rho(x,y)} L(d_\rho, \alpha) = -\beta \left( \log \frac{d_\rho(x, y)}{d_\rho^\mu(x, y)} - \mathbf{1} \right) + e_\alpha(x, y),$$

which, when set to zero, yields

$$d_\rho(x, y) = d_\rho^\mu(x, y) \exp \left( \frac{1}{\beta} e_\alpha(x, y) \right) \exp(-1).$$

Primal feasibility implies that  $\sum_{x,y} d_\rho(x, y) = 1$ , thus, our choice of  $\alpha$  should satisfy such condition. Letting  $Z = \exp(1)$ , and  $\alpha^*$  be the optimal Lagrange multiplier, we have

$$d_\rho^*(x, y) = \frac{1}{Z} d_\rho^\mu(x, y) \exp \left( \frac{1}{\beta} e_{\alpha^*}(x, y) \right). \quad (13)$$

Writing the expression for the reward function, we get

$$r^*(x, y) = \beta \log \frac{d_\rho^*(x, y)}{d_\rho^\mu(x, y)} + \beta + \gamma \sum_{x'} \mathbf{1}(x = T(x', y')) \alpha^*(x') - \alpha^*(x).$$

Now, observe that, given a trajectory  $\tau = (x_0, y_0, x_1, \dots)$ , we can write the discounted return using the above expression and obtain

$$\begin{aligned} \sum_{t=0}^{\infty} \gamma^t r^*(x_t, y_t) &= \sum_{t=0}^{\infty} \gamma^t \left( \beta \log \frac{d_\rho^*(x_t, y_t)}{d_\rho^\mu(x_t, y_t)} + \beta + \gamma \sum_x \mathbf{1}(x = T(x', y')) \alpha^*(x) - \alpha^*(x_t) \right) \\ &= \sum_{t=0}^{\infty} \gamma^t \left( \beta \log \frac{d_\rho^*(x_t, y_t)}{d_\rho^\mu(x_t, y_t)} + \beta + \gamma \alpha^*(x_{t+1}) - \alpha^*(x_t) \right) \end{aligned} \quad (14)$$

$$= \sum_{t=0}^{\infty} \gamma^t \left( \beta \log \frac{d_\rho^*(x_t, y_t)}{d_\rho^\mu(x_t, y_t)} + \beta - \alpha^*(x_0) \right), \quad (15)$$

where for Equation (14) we have used the fact that the transitions are deterministic, and for Equation (15) note that the terms  $\alpha^*(x)$  cancel each other out.

Now, let us get back to the BT preference model for MDPs. Given a dataset  $\mathcal{D}_n$  of pairs of trajectories, each pair of which starts from the same initial state, we can express the MLE loss directly in terms of the occupancy measures using the above derivation as follows:

$$\mathcal{L}_{\text{DPO}}(d_\rho) = -\mathbb{E}_{(\tau^w, \tau^l) \sim \mathcal{D}_n} \left[ \log \sigma \left( \sum_{t=0}^{\infty} \gamma^t \beta \log \frac{d_\rho(x_t^w, y_t^w)}{d_\rho^\mu(x_t^w, y_t^w)} - \sum_{t=0}^{\infty} \gamma^t \beta \log \frac{d_\rho(x_t^l, y_t^l)}{d_\rho^\mu(x_t^l, y_t^l)} \right) \right],$$

where we have used the fact that the terms  $\beta$  and  $\alpha^*(x_0)$  cancel out.

Now, note that the minimizer to the above loss may not satisfy the Bellman flow constraints of Problem (P3.2'). Thus, we need to restrict the domain of the problem to the following set

$$\mathcal{B} = \left\{ d \in \Delta(\mathcal{X} \times \mathcal{Y}) : \sum_y d(x, y) = (1 - \gamma)\rho(x) + \gamma \sum_{x', y'} \mathbf{1}(x = T(x', y')) d(x', y'), \forall x \in \mathcal{X} \right\}$$

### H.3. DPO for MDPs with Loglinear Occupancy Measures

Similar to the contextual bandit setting, we want to write the DPO loss such that it resembles logistic regression. For loglinear occupancy measures, as defined in Definition 7.1, with parameter set restricted to

$$\Theta' := \{ \theta \in \mathbb{R}^{d_M} : d_\rho^{\pi_\theta} \in \mathcal{B} \},$$

we can write the loss so that it resembles logistic regression. Note that the domain of  $\theta$  is restricted only to those parameters which imply that  $d_\rho^{\pi_\theta}$  is an occupancy measure with respect to the underlying MDP. For this case, the DPO loss becomes

$$\mathcal{L}_{\mathcal{D}_n}(\theta) = -\mathbb{E}_{(\tau^w, \tau^l) \sim \mathcal{D}_n} \left[ \log \sigma \left( \beta \theta^\top \left( \sum_{t=0}^{\infty} \gamma^t (\psi(x_t^w, y_t^w) - \psi(x_t^l, y_t^l)) \right) + K(\tau^w, \tau^l) \right) \right]$$

where

$$K(\tau^w, \tau^l) = \sum_{t=0}^{\infty} \gamma^t \log \frac{d_\rho^\mu(x_t^l, y_t^l)}{d_\rho^\mu(x_t^w, y_t^w)}.$$

Given the learned occupancy measure  $d_\rho^{\pi_\theta}$ , one can finally compute an optimal policy, for each state-action pair, as

$$\pi_\theta(y|x) = \frac{d_\rho^{\pi_\theta}(x, y)}{\sum_y d_\rho^{\pi_\theta}(x, y)}.$$

Note that, in general, the quantity  $K(\tau^w, \tau^l)$  is not easy to compute as it requires access to the occupancy measure with respect to  $\mu$ . However, in practice,  $K(\tau^w, \tau^l)$  can be treated as a hyperparameter of the problem and tuned accordingly.

## I. Gradient Expression for KL-regularized Objective in MDPs

In this section, we derive the gradient for the loglinear policy class. We rewrite the problem below for convenience.

$$\max_{\theta} \mathbb{E}_{x \sim \rho} \left[ \sum_{t=0}^{\infty} \gamma^t (r(x_t, y_t) - \beta D_{\text{KL}}(\pi_\theta(\cdot|x_t) \parallel \mu(\cdot|x_t))) \mid y_t \sim \pi_\theta(\cdot|x_t) \right]$$

**Lemma I.1.** *Let*

$$\mathcal{V}_r^{\pi_\theta}(x) = \mathbb{E}_{x \sim \rho, y_t \sim \pi_\theta(\cdot|x_t)} \left[ \sum_{t \geq 0} \gamma^t \left( r(x_t, y_t) - \beta \log \frac{\pi_\theta(y_t|x_t)}{\mu(y_t|x_t)} \right) \right]$$

and

$$\mathcal{Q}_r^{\pi_\theta}(x, y) = r(x, y) + \gamma \mathcal{V}_r^{\pi_\theta}(T(x, y)).$$

The gradient expression for  $\mathcal{V}^{\pi_\theta}(\rho)$  is given by

$$\nabla_{\theta} \mathcal{V}_r^{\pi_\theta}(\rho) = \frac{1}{1-\gamma} \sum_x d_\rho^{\pi_\theta}(x) \sum_y \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \bar{\psi}_\theta(x, y).$$

*Proof.* Note that we have

$$\mathcal{V}_r^{\pi_\theta}(\rho) = \mathbb{E}_{x \sim \rho} \left[ \sum_y \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right].$$

Thus, we can write

$$\begin{aligned} \nabla_{\theta} \mathcal{V}_r^{\pi_\theta}(\rho) &= \sum_{x, y} \rho(x) \left( \nabla_{\theta} \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \right. \\ &\quad \left. + \pi_\theta(y|x) \left( \nabla_{\theta} \mathcal{Q}_r^{\pi_\theta}(x, y) - \frac{\mu(y|x)}{\pi_\theta(y|x)} \nabla_{\theta} \pi_\theta(y|x) \right) \right) \\ &= \sum_{x, y} \rho(x) \left( \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x, y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} - 1 \right) \bar{\psi}_\theta(x, y) + \pi_\theta(y|x) \nabla_{\theta} \mathcal{Q}_r^{\pi_\theta}(x, y) \right) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{x,y} \rho(x) \left( \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x,y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \bar{\psi}_\theta(x,y) \right) \\
 &\quad + \gamma \sum_{x,y} \rho(x) \pi_\theta(y|x) \nabla_\theta \mathcal{V}_r^{\pi_\theta}(T(x,y)) \\
 &= \frac{1}{1-\gamma} \sum_x d_\rho^{\pi_\theta}(x) \sum_y \pi_\theta(y|x) \left( \mathcal{Q}_r^{\pi_\theta}(x,y) - \beta \log \frac{\pi_\theta(y|x)}{\mu(y|x)} \right) \bar{\psi}_\theta(x,y),
 \end{aligned}$$

where the second equality follows from the derivation of the gradient of loglinear policies (see the proof of Lemma E.1, while the third equality follows from the fact that  $\mathbb{E}_{y \sim \pi_\theta(\cdot|x)}[\bar{\psi}_\theta(x,y)] = 0$ , for each  $x \in \mathcal{X}$ .  $\square$

## J. Technical Lemmas

The purpose of this section is to present various technical results that are useful for our paper. Let us denote by  $\Phi \in \mathbb{R}^{d_R \times XY}$  and  $\Psi \in \mathbb{R}^{d_P \times XY}$  the reward and policy feature matrices with columns  $\phi(x,y)$  and  $\psi(x,y)$ , respectively.

**Lemma J.1.** *Assume that  $r^* \in \mathcal{F}$ ,  $\pi^* \in \Pi$  and  $\mu \in \Pi$ , for some  $\pi_{r^*}^* \in \arg \max_\pi \mathcal{V}^\pi(\rho)$ . Furthermore, assume that the column space of  $\Phi$  is a subspace of the column space of  $\Psi$ . Then, there exists  $\theta^* \in \Theta$ , for which  $\pi_{\theta^*}$  maximizes the objective of (P1.2) and that can be represented in terms of the ground-truth reward function, i.e.  $\pi_{r^*}^*(y|x) = \pi_{\theta^*}(y|x) \propto \mu(y|x) \exp(r^*(x,y)/\beta)$ , for all  $(x,y)$ .*

*Proof.* Let  $x \in \mathcal{X}$ . From Equation (1) we have that

$$\begin{aligned}
 \pi_{r^*}^*(x,y) &\propto \mu(y|x) \exp\left(\frac{1}{\beta} r^*(x,y)\right) \\
 &\propto \exp\left(\theta_\mu^\top \psi(x,y) + \phi(x,y)^\top \omega^*\right),
 \end{aligned}$$

for some  $\pi_{r^*}^* \in \arg \max_\pi \mathcal{V}_r^\pi(\rho)$ , where the second relation holds due to the assumptions on the policy class and reward class. Thus, if we can find a  $\theta^* \in \mathbb{R}^{d_P}$  such that, for all  $(x,y)$ , we have

$$\exp\left(\psi(x,y)^\top \theta^*\right) = \exp\left(\theta_\mu^\top \psi(x,y) + \phi(x,y)^\top \omega^*\right),$$

then we have shown that the optimal regularized policy belongs to the loglinear class. For the above to hold, we equivalently need

$$\Psi^\top (\theta^* - \theta_\mu) - \Phi^\top \omega^* = \mathbf{0}.$$

The above equation has a solution for  $\theta^*$  whenever the column space of  $\Phi$  is contained in the column space of  $\Psi$ .  $\square$

**Lemma J.2.** *Assume that  $r^* \in \mathcal{F}$ ,  $d_\rho^\mu \in \Pi'$  and  $d_{\rho^*}^{\pi_{r^*}^*} \in \Pi'$ , for some optimal  $d_{\rho^*}^{\pi_{r^*}^*}$ . Furthermore, assume that the column space of  $\Phi + \Phi_{\pi_{r^*}^*}$  is contained in the column space of  $\Psi$ , where  $\Phi_{\pi_{r^*}^*} \in \mathbb{R}^{d_R \times XY}$  has columns*

$$\gamma \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x_t, y_t) \middle| x_0 = x, \pi_{r^*}^* \right] - \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x_t, y_t) \middle| x_0 = x, y_0 = T(x,y), \pi_{r^*}^* \right].$$

*Then, for finite MDPs with deterministic transitions, there exists  $\theta^*$  such that  $d_{\rho^*}^{\pi_{r^*}^*}(x,y)$ .*

*Proof.* From Equation (13) we have

$$d_\rho^*(x,y) = \frac{1}{Z} d_\rho^\mu(x,y) \exp\left(\frac{1}{\beta} e_{\alpha^*}(x,y)\right),$$

where

$$e_\alpha(x,y) = r^*(x,y) + \gamma \alpha^*(T(x,y)) - \alpha^*(x),$$

is the advantage function when using  $\alpha^*$  and  $\alpha^*$  denote the optimal dual variables for Problem (P3.2'). As shown in (Lee et al., 2021), these variables correspond to the optimal value function with respect to  $r^*$ . Thus, we have

$$\begin{aligned}
 e_{\alpha^*}(x, y) &= r^*(x, y) + \gamma \alpha^*(T(x, y)) - \alpha^*(x) \\
 &= \phi(x, y)^\top \omega^* + \gamma \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r^*(x, y) \middle| x_0 = x, \pi_{r^*}^* \right] - \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r^*(x, y) \middle| x_0 = x, y_0 = T(x, y), \pi_{r^*}^* \right] \\
 &= \phi(x, y)^\top \omega^* + \gamma \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y)^\top \omega^* \middle| x_0 = x, \pi_{r^*}^* \right] - \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y)^\top \omega^* \middle| x_0 = x, y_0 = T(x, y), \pi_{r^*}^* \right] \\
 &= \left( \phi(x, y) - \gamma \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y) \middle| x_0 = x, \pi_{r^*}^* \right] - \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y) \middle| x_0 = x, y_0 = T(x, y), \pi_{r^*}^* \right] \right)^\top \omega^* .
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 d_\rho^*(x, y) &= \frac{1}{Z} d_\rho^\mu(x, y) \exp \left( \frac{1}{\beta} e_{\alpha^*}(x, y) \right) \propto \exp \left( \theta_\mu^\top \psi(x, y) \right. \\
 &\quad \left. + (\omega^*)^\top \left( \phi(x, y) + \gamma \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y) \middle| x_0 = x, \pi_{r^*}^* \right] - \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t \phi(x, y) \middle| x_0 = x, y_0 = T(x, y), \pi_{r^*}^* \right] \right) \right) .
 \end{aligned}$$

For the above to hold, we equivalently need

$$\Psi^\top (\theta^* - \theta_\mu) - \left( \Phi + \Phi_{\pi_{r^*}^*} \right)^\top \omega^* = \mathbf{0} .$$

The above has a solution whenever the column space of  $\Phi + \Phi_{\pi_{r^*}^*}$  is contained in the column space of  $\Psi$ .  $\square$

Next, we will prove a result that connects the suboptimality gap with the gap in terms of the KL-regularized objectives.

**Lemma J.3.** *For any  $\theta$ , we have*

$$\beta D_{\text{KL}}(\pi_{\theta^*} \| \mu) - \beta D_{\text{KL}}(\pi_\theta \| \mu) \leq D(\pi_\theta) \leq \beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \| \mu) - \beta D_{\text{KL}}(\pi_\theta \| \mu) ,$$

*Proof.* Note that, by definition,

$$\begin{aligned}
 G(\pi_\theta) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_\theta}(\rho) \\
 &= (V_{r^*}^{\text{opt}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho)) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) + (\mathcal{V}_{r^*}^{\pi_\theta}(\rho) - V_{r^*}^{\pi_\theta}(\rho)) \\
 &\leq \left( V_{r^*}^{\text{opt}}(\rho) - \mathcal{V}_{r^*}^{\pi_{r^*}^{\text{opt}}}(\rho) \right) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) + (\mathcal{V}_{r^*}^{\pi_\theta}(\rho) - V_{r^*}^{\pi_\theta}(\rho)) \\
 &\leq \beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \| \mu) - \beta D_{\text{KL}}(\pi_\theta \| \mu) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) .
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 G(\pi_\theta) &= V_{r^*}^{\text{opt}}(\rho) - V_{r^*}^{\pi_\theta}(\rho) \\
 &= (V_{r^*}^{\text{opt}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho)) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) + (\mathcal{V}_{r^*}^{\pi_\theta}(\rho) - V_{r^*}^{\pi_\theta}(\rho)) \\
 &\geq (V_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho)) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) + (\mathcal{V}_{r^*}^{\pi_\theta}(\rho) - V_{r^*}^{\pi_\theta}(\rho)) \\
 &\geq \beta D_{\text{KL}}(\pi_{\theta^*} \| \mu) - \beta D_{\text{KL}}(\pi_\theta \| \mu) + (\mathcal{V}_{r^*}^{\pi_{\theta^*}}(\rho) - \mathcal{V}_{r^*}^{\pi_\theta}(\rho)) .
 \end{aligned}$$

The result follows.  $\square$

Next, we will control the quantity  $D(\pi_\theta)$  for the DPO setting.

**Lemma J.4.** *With probability at least  $1 - \delta$ , we have*

$$\begin{aligned} D(\pi_{\hat{\theta}}) &\leq \beta D_{\text{KL}}(\pi_{r^*}^{\text{opt}}, \pi_{\theta^*}) + \frac{\Lambda_P U d_P}{S_P n} \sqrt{\frac{\log(4/\delta)}{2n}} \\ &= \beta \left( D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \|\mu) - D_{\text{KL}}(\pi_{\theta^*} \|\mu) \right) + O\left(\frac{d_P}{n^{3/2}}\right). \end{aligned}$$

*Proof.* Recall that, for any  $\theta$ , we have defined

$$D_{\text{KL}}(\pi_{\theta} \|\mu) = \sum_x \rho(x) \sum_{x,y} \pi_{\theta}(y|x) \log \frac{\pi_{\theta}(y|x)}{\mu(y|x)}.$$

First, note that

$$\begin{aligned} D_{\text{KL}}(\pi_{\theta^*} \|\mu) - D_{\text{KL}}(\pi_{\theta} \|\mu) &= \left( D_{\text{KL}}(\pi_{\theta^*} \|\mu) - \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \right) + \left( \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \right. \\ &\quad \left. - \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\hat{\theta}} \|\mu) \right) + \left( \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\hat{\theta}} \|\mu) - D_{\text{KL}}(\pi_{\theta} \|\mu) \right) \\ &= \left( D_{\text{KL}}(\pi_{\theta^*} \|\mu) - \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \right) + \left( \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\hat{\theta}} \|\mu) - D_{\text{KL}}(\pi_{\theta} \|\mu) \right), \end{aligned}$$

where the second equality follows from the exact optimization assumption. Note that the two summands above are deviations from means. If we can show that each individual quantity is bounded, then we can apply Hoeffding bounds. To that end, first, note that

$$\begin{aligned} D_{\text{KL}}(\pi_{\theta^*} \|\mu) &= \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} = \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \log \frac{\mu(y|x) \exp\left(\frac{1}{\beta} r^*(x, y)\right)}{\mu(y|x)} \\ &= \frac{1}{\beta} \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) r^*(x, y) \leq \frac{1}{\beta}, \end{aligned}$$

since the reward cannot be more than 1. Similarly, for every  $x \in \mathcal{D}$ , we have

$$\begin{aligned} D_{\text{KL}}(\pi_{\hat{\theta}} \|\mu) &= D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) = \sum_x \rho(x) \sum_y \pi_{\theta_{\mathcal{D}_n}^*}(y|x) \log \frac{\mu(y|x) \exp\left(\frac{1}{\beta} r^*(x, y)\right)}{\mu(y|x)} \\ &= \frac{1}{\beta} \sum_x \rho(x) \sum_y \pi_{\theta_{\mathcal{D}_n}^*}(y|x) r^*(x, y) \leq \frac{1}{\beta}. \end{aligned}$$

For other contexts  $x \notin \mathcal{D}$ , such a relation does not hold. Thus, we take another approach. We show that, for such points, the KL divergence between the learned policy and the sampling policy cannot be too far away from that between the optimal policy and the sampling policy.

$$\begin{aligned} \left| D_{\text{KL}}(\pi_{\theta^*} \|\mu) - D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \right| &= \left| \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \log \frac{\pi_{\theta^*}(y|x)}{\mu(y|x)} - \sum_{x,y} \pi_{\theta_{\mathcal{D}_n}^*}(y|x) \log \frac{\pi_{\theta_{\mathcal{D}_n}^*}(y|x)}{\mu(y|x)} \right| \\ &= \left| \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \left( \log \frac{\pi_{\theta^*}(y|x)}{\pi_{\theta_{\mathcal{D}_n}^*}(y|x)} + \log \exp\left(\frac{1}{\beta} r^*(x, y)\right) \right) \right. \\ &\quad \left. - \sum_{x,y} \pi_{\theta_{\mathcal{D}_n}^*}(y|x) \left( \log \frac{\pi_{\theta_{\mathcal{D}_n}^*}(y|x)}{\pi_{\theta^*}(y|x)} + \log \exp\left(\frac{1}{\beta} r^*(x, y)\right) \right) \right| \\ &\leq \frac{1}{\beta} \left| \left( V_{r^*}^{\pi_{\theta^*}}(\rho) - V^{\pi_{\theta_{\mathcal{D}_n}^*}}(\rho) \right) - D_{\text{KL}}(\pi_{\theta^*} \|\pi_{\theta_{\mathcal{D}_n}^*}) \right| \end{aligned}$$

$$\leq \frac{1}{\beta} + D_{\text{KL}}(\pi_{\theta^*} \|\pi_{\theta_{\mathcal{D}_n}^*}) .$$

Now, for the last term of the right-hand side, we have

$$\begin{aligned} D_{\text{KL}}(\pi_{\theta^*} \|\pi_{\theta_{\mathcal{D}_n}^*}) &= \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \left( \log \pi_{\theta^*}(y|x) - \log \pi_{\theta_{\mathcal{D}_n}^*}(y|x) \right) \\ &= \sum_x \rho(x) \sum_y \pi_{\theta^*}(y|x) \left( \langle \psi(x, y), \theta^* - \theta_{\mathcal{D}_n}^* \rangle + \log \frac{\sum_{x', y'} \exp(\psi(x', y')^\top \theta_{\mathcal{D}_n}^*)}{\sum_{x', y'} \exp(\psi(x', y')^\top \theta^*)} \right) \\ &\leq \frac{\Lambda_P U d_P}{S_P \beta n} , \end{aligned}$$

where the last inequality follows from the same arguments as in the proof of Theorem 4.2. Going back to the original expression, note that, for any given  $x \in \mathcal{X}$ , we have

$$0 \leq D_{\text{KL}}(\pi_{\theta^*} \|\mu) \leq \frac{1}{\beta} , \quad \text{and} \quad 0 \leq D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \leq \frac{1}{\beta} + \frac{\Lambda_P U d_P}{S_P \beta n} .$$

Thus, by Hoeffding's inequality, for any  $\delta \geq 0$ , with probability at least  $1 - \delta$ , we have

$$\left| D_{\text{KL}}(\pi_{\theta^*} \|\mu) - \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\theta_{\mathcal{D}_n}^*} \|\mu) \right| \leq \frac{1}{\beta} \sqrt{\frac{\log(4/\delta)}{2n}} ,$$

and

$$\left| \frac{1}{n} \sum_{x \in \mathcal{D}} D_{\text{KL}}(\pi_{\tilde{\theta}} \|\mu) - D_{\text{KL}}(\pi_{\theta} \|\mu) \right| \leq \left( \frac{1}{\beta} + \frac{\Lambda_P U d_P}{S_P \beta n} \right) \sqrt{\frac{\log(4/\delta)}{2n}} ,$$

which implies that

$$- \left( \frac{2}{\beta} + \frac{\Lambda_P U d_P}{S_P \beta n} \right) \sqrt{\frac{\log(4/\delta)}{2n}} \leq D_{\text{KL}}(\pi_{\theta^*} \|\mu) - D_{\text{KL}}(\pi_{\theta} \|\mu) \leq \left( \frac{2}{\beta} + \frac{\Lambda_P U d_P}{S_P \beta n} \right) \sqrt{\frac{\log(4/\delta)}{2n}} .$$

On the other hand, note that

$$\begin{aligned} D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \|\mu) - D_{\text{KL}}(\pi_{\tilde{\theta}} \|\mu) &= (D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \|\mu) - D_{\text{KL}}(\pi_{\theta^*} \|\mu)) + (D_{\text{KL}}(\pi_{\theta^*} \|\mu) - D_{\text{KL}}(\pi_{\tilde{\theta}} \|\mu)) \\ &\leq (D_{\text{KL}}(\pi_{r^*}^{\text{opt}} \|\mu) - D_{\text{KL}}(\pi_{\theta^*} \|\mu)) + \left( \frac{2}{\beta} + \frac{\Lambda_P U d_P}{S_P \beta n} \right) \sqrt{\frac{\log(4/\delta)}{2n}} . \end{aligned}$$

□

Next, we prove some useful properties of the log-exp-sum function.

**Lemma J.5.** *The function defined as*

$$A(\theta) = \sum_x \rho(x) \log \sum_{x, y} \exp(\theta^\top \psi(x, y)) .$$

*is 1-Lipschitz and 2-smooth. Moreover, if the features are sampled from a 0-mean distribution and span  $\mathbb{R}^{d_P}$ , then there exists  $\kappa > 0$ , such that  $A(\theta)$  is  $\kappa$ -strongly convex.*

*Proof.* Let  $\theta \in \mathbb{R}^{d_P}$ . Note that

$$\nabla_{\theta} A(\theta) = \sum_x \rho(x) \frac{\sum_y \exp(\psi(x, y)^\top \theta)}{\sum_{y'} \exp(\psi(x, y')^\top \theta)} \psi(x, y')$$



$$\begin{aligned}
 &= \sum_x \rho(x) \sum_y \pi_\theta(y|x) \psi(x, y) \\
 &\leq \max_{x, y} \|\psi(x, y)\|_2 \\
 &\leq 1.
 \end{aligned}$$

On the other hand, the Hessian of  $A(\theta)$  is

$$\begin{aligned}
 \nabla_\theta^2 A(\theta) &= \sum_x \rho(x) \sum_y \nabla_\theta \pi_\theta(y|x) \psi(x, y) \\
 &= \sum_x \rho(x) \sum_y \pi_\theta(y|x) (\psi(x, y) - \mathbb{E}_{y' \sim \pi_\theta(\cdot|x)}[\psi(x, y')]) \psi(x, y)^\top \\
 &= \mathbb{E}_{x \sim \rho, y \sim \pi_\theta(\cdot|x)} [\psi(x, y) \psi(x, y)^\top] - \mathbb{E}_{x \sim \rho, y \sim \pi_\theta(\cdot|x)} [\psi(x, y)] \mathbb{E}_{x \sim \rho, y \sim \pi_\theta(\cdot|x)} [\psi(x, y)]^\top \\
 &= \mathbb{E}_{x \sim \rho, y \sim \pi_\theta(y|x)} \left[ (\psi(x, y) - \mathbb{E}_\theta [\psi(x, y)]) (\psi(x, y) - \mathbb{E}_\theta [\psi(x, y)])^\top \right].
 \end{aligned}$$

By assumption on the feature mapping, we have that

$$\begin{aligned}
 \|\nabla_\theta^2 A(\theta)\|_2 &\leq \max_{x, y} \left\| (\psi(x, y) - \mathbb{E}_\theta [\psi(x, y)]) (\psi(x, y) - \mathbb{E}_\theta [\psi(x, y)])^\top \right\|_2 \\
 &\leq \max_{x, y} \|\psi(x, y) - \mathbb{E}_\theta [\psi(x, y)]\|_2 \\
 &\leq 2 \max_{x, y} \|\psi(x, y)\|_2 = 2.
 \end{aligned}$$

Therefore, the function  $A(\theta)$  is 2-smooth in  $\theta$ . For strong convexity, let  $\psi$  be sampled from a 0-mean bounded distribution. Note that, for any non-zero vector in  $\mathbb{R}^{d_P}$ , we have

$$\begin{aligned}
 z^\top \nabla_\theta^2 A(\theta) z &= \mathbb{E}_{x \sim \rho, y \sim \pi_\theta(\cdot|x)} [z^\top \psi(x, y) \psi(x, y)^\top z] \\
 &\geq \min_{\theta, x, y} \pi_\theta(\cdot|x) \sum_{x, y} (\psi(x, y)^\top z)^2 \\
 &\geq C_3 \sum_{x, y} (\psi(x, y)^\top z)^2,
 \end{aligned}$$

for a positive  $C_3$ , since  $\pi_\theta$  is in the loglinear class, for every  $\theta$ , and using Lemma E.7. Now, note that, if  $z$  can be expressed as a linear combination of  $\{\psi(x, y)\}_{x, y}$ , the summation cannot be zero for non-zero  $z$ . Thus, if  $\{\psi(x, y)\}_{x, y}$  spans  $\mathbb{R}^{d_P}$ , that is, the feature matrix is full rank, then there exists an absolute positive constant  $\kappa$ , such that we have

$$\|\nabla_\theta^2 A(\theta)\|_2 \geq \kappa > 0.$$

Thus, the function  $A(\theta)$  is  $\kappa$ -strongly convex.  $\square$

**Lemma J.6.** *In general, the norms of the gradient and Hessian for the loss of tabular DPO are unbounded from above.*

*Proof.* Observe that, given policy  $\pi$  and  $(x, y^w) \in \mathcal{D}$ , we have

$$\nabla_{\pi(y^w|x)} \mathcal{L}_{\mathcal{D}}(\pi) = \frac{\beta}{n} \left( 1 - \sigma \left( \beta \log \frac{\pi(y^w|x)}{\mu(y^w|x)} - \beta \log \frac{\pi(y^l|x)}{\mu(y^l|x)} \right) \right) \frac{1}{\pi(y^w|x)}.$$

On the other hand, for the second derivative with respect to  $\pi(y^w|x)$ , we have the following. First, let

$$f(\pi(y^w|x)) = \beta \log \frac{\pi(y^w|x)}{\mu(y^w|x)} - \beta \log \frac{\pi(y^l|x)}{\mu(y^l|x)}.$$

We have that  $\nabla_{\pi(y^w|x)} f(\pi(y^w|x)) = \beta/\pi(y^w|x)$ . Now, observe that

$$\nabla_{\pi(y^w|x)}^2 \mathcal{L}_{\mathcal{D}}(\pi) = \frac{\beta}{n} \nabla_{\pi(y^w|x)} \frac{\exp(f(\pi(y^w|x)))}{\pi(y^w|x)(1 + \exp(f(\pi(y^w|x))))}$$

$$\begin{aligned}
 &= \frac{\beta}{n} \left( \frac{\frac{\beta}{\pi(y^w|x)} \exp(f(\pi(y^w|x))) \pi(y^w|x)(1 + \exp(f(\pi(y^w|x))))}{(\pi(y^w|x)(1 + \exp(f(\pi(y^w|x))))^2)} \right) \\
 &\quad - \frac{\beta}{n} \left( \frac{\exp(f(\pi(y^w|x))) \left( (1 + \exp(f(\pi(y^w|x)))) + \pi(y^w|x) \frac{\beta}{\pi(y^w|x)} \exp(f(\pi(y^w|x))) \right)}{(\pi(y^w|x)(1 + \exp(f(\pi(y^w|x))))^2)} \right) \\
 &= \frac{\beta \left( (\beta - 1) \exp(f(\pi(y^w|x)))(1 + \exp(f(\pi(y^w|x)))) - \beta \exp(f(\pi(y^w|x)))^2 \right)}{n (\pi(y^w|x)(1 + \exp(f(\pi(y^w|x))))^2)}.
 \end{aligned}$$

The above numerator is not always non-negative, as solving for  $\exp(f(\pi(y^w|x)))$  will show. Moreover, neither the norm of the gradient nor the operator norm of the Hessian can be upper-bounded in general, due to the presence of  $\pi(y^w|x)$  in the denominator.  $\square$

**Theorem J.7** (Theorem 1.(c) of (Shah et al., 2016)). *For the BT preference model,  $B$ -bounded weight vector and sample size  $n \geq O(\text{tr}(\Sigma^\dagger)/\beta^2 B^2)$ , where  $\Sigma$  denotes the Laplacian with respect to features, the maximum likelihood estimator satisfies the minimax bounds*

$$\Omega\left(\frac{d}{\beta n}\right) \leq \|\tilde{\theta} - \theta^*\|_{\Sigma}^2 \leq O\left(\frac{d}{\beta n}\right).$$

**Lemma J.8** (Lemma 27 of (Mei et al., 2020)). *Let  $\pi_{\theta} = \text{softmax}(\Psi\theta)$  and  $\pi_{\theta'} = \text{softmax}(\Psi\theta')$ . Then, for any constant  $c$ , we have*

$$D_{\text{KL}}(\pi_{\theta}||\pi_{\theta'}) \leq \frac{1}{2} \|\Psi\theta - \Psi\theta' - c^{\top} \mathbf{1}\|^2.$$