

# SHARP ANALYSIS FOR KL-REGULARIZED CONTEXTUAL BANDITS AND RLHF

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

*Reverse-Kullback-Leibler* regularization has emerged to be a predominant technique used to enhance policy optimization in reinforcement learning (RL) and reinforcement learning from human feedback (RLHF), which forces the learned policy to stay close to a reference policy. While the effectiveness and necessity of KL-regularization has been empirically demonstrated in various practical scenarios, current theoretical analysis of KL-regularized RLHF still obtain the same  $\mathcal{O}(1/\epsilon^2)$  sample complexity as problems without KL-regularization. To understand the fundamental distinction between policy learning objectives with KL-regularization and ones without KL-regularization, we are the first to theoretically demonstrate the power of KL-regularization by providing a sharp analysis for KL-regularized contextual bandits and RLHF, revealing an  $\mathcal{O}(1/\epsilon)$  sample complexity when  $\epsilon$  is sufficiently small.

We further explore the role of data coverage in contextual bandits and RLHF. While the coverage assumption is commonly employed in offline RLHF to link the samples from the reference policy to the optimal policy, often at the cost of a multiplicative dependence on the coverage coefficient, its impact on the sample complexity of online RLHF remains unclear. Previous theoretical analyses of online RLHF typically require explicit exploration and additional structural assumptions on the reward function class. In contrast, we show that with sufficient coverage from the reference policy, a simple two-stage mixed sampling strategy can achieve a sample complexity with only an additive dependence on the coverage coefficient. Our results provide a comprehensive understanding of the roles of KL-regularization and data coverage in RLHF, shedding light on the design of more efficient RLHF algorithms.

## 1 INTRODUCTION

Recently, *Reinforcement Learning from Human Feedback* (RLHF) has emerged as a central tool for aligning large language models (LLMs) and diffusion models with human values and preferences (Christiano et al., 2017; Ziegler et al., 2019; Ouyang et al., 2022; Bai et al., 2022; Rafailov et al., 2024), exhibiting impressive capabilities in applications, such as Chatgpt (Achiam et al., 2023), Claude (Anthropic, 2023), Gemini (Team et al., 2023), and LLaMA-3 (Meta, 2024).

RLHF methods treat the language model as a policy that takes a prompt  $x$  and produces a response  $a$  conditioned on  $x$ , and they optimize the policy by aligning it with human feedback. There are two main kinds of feedback: absolute rating and preference comparison. In practice, collecting the absolute ratings typically involving the human annotators to provide rating scores like 1 to 5 (Wang et al., 2024a;b) for the responses or hard 0-1 scores for math reasoning tasks since the reasoning tasks often have golden answers (Cobbe et al., 2021; Hendrycks et al., 2021; Xiong et al., 2024b). Additionally, preference comparison is frequently applied in chat tasks when making comparisons is much easier for human labler (Achiam et al., 2023).

Since the human value and preference are so complicated that they are unlikely to be encompassed by the considered preference model classes (such as the absolute reward model or the reward-based *Bradley-Terry* model (Bradley and Terry, 1952a)), the learned model is easy to be hacked and biased. Practically, the policy may generate disproportionate bold words or emoji to please the learned reward (Zhang et al., 2024). Hence, the KL-regularization between the learned policy and a reference policy (the pre-trained model after supervised fine-tuning) plays a fundamental role in RLHF to

054 avoid overfitting. There is a line of RLHF work that realizes the significance of KL-regularization  
 055 and regards the problem as a reverse-KL regularized contextual bandit (Ziegler et al., 2019; Wu et al.,  
 056 2021; Ouyang et al., 2022; Rafailov et al., 2024; Xiong et al., 2024a; Ye et al., 2024b). However,  
 057 they basically adopt the techniques from bandit framework and neglect the characteristic of reverse-  
 058 KL-regularization, thus obtaining almost the same sample complexity with problems without KL-  
 059 regularization. Therefore, the question of *whether there exists a fundamental distinction between*  
 060 *policy learning objectives with KL-regularization and ones without KL-regularization* is still largely  
 061 under-explored.

062 Compared to the offline RLHF algorithms (Rafailov et al., 2024; Azar et al., 2024; Chen et al.,  
 063 2024) that can only use planning to approximate the solution of the minimizing relative entropy  
 064 optimization (Ziebart et al., 2008; Song et al., 2024), online RLHF has been demonstrated to out-  
 065 perform offline methods empirically and theoretically (Bai et al., 2023; Meta, 2024; Xiong et al.,  
 066 2024a; Tajwar et al., 2024; Song et al., 2024), because it has further interactions with human or the  
 067 preference oracle. Most standard theoretical online RL techniques apply optimism in the balance  
 068 of exploration and exploitation (Abbasi-Yadkori et al., 2011; Wang et al., 2020). However, it is  
 069 inefficient to implement exploration for practical RLHF algorithms. Meanwhile, an emerging line  
 070 of offline RLHF literature highlights the coverage of the reference policy  $\pi_0$ . The coverage of  $\pi_0$   
 071 refers to the ability of the model to generate diverse responses for a wide range of prompts. A model  
 072 with good coverage can generalize well to unseen contexts and actions, which is essential for the  
 073 learned reward function to generalize well. In practice, this is evidenced by the fact that the simple  
 074 best-of-n sampling based on  $\pi_0$  is competitive with the well-tuned PPO algorithm for general  
 075 open-ended conversation tasks (Dong et al., 2023), and the fact that the  $\pi_0$  can solve a majority of  
 076 the math problems with multiple responses (Shao et al., 2024; Nakano et al., 2021). However, the  
 077 theoretical understanding of the role of coverage in online RLHF is still largely understudied. Thus,  
 078 it is natural to ask *is explicit exploration necessary for online RLHF with a good coverage of  $\pi_0$  and*  
 079 *how the coverage of  $\pi_0$  affects the sample complexity of online RLHF.*

079 In this paper, we answer the above questions by

- 081 • providing a novel fine-grained analysis for KL-regularized in contextual bandits and RLHF, which  
 082 adapts to the optimization landscape of the reverse-KL regularization and reveals a sharper sample  
 083 complexity than the existing results, and
- 084 • proposing an efficient 2-stage mixed sampling strategy for online RLHF with a good coverage  
 085 of  $\pi_0$ , which achieves a sample complexity with only an additive dependence on the coverage  
 086 coefficient.

### 087 1.1 OUR CONTRIBUTIONS

088 In this paper, we make a first attempt to illustrate the statistical benefits of KL-regularization for  
 089 policy optimization in contextual bandits and reinforcement learning from preference feedback.

090 Our main contributions are summarized as follows:

- 091 • In Section 3, we formulate RLHF with absolute-rating feedback as a contextual bandit problem  
 092 with KL-regularization. First, we provide a novel lower bound for the KL-regularized contextual  
 093 bandit problem, which indicates that the sample complexity of the problem is  $\Omega(\eta \log N_{\mathcal{R}}(\epsilon)/\epsilon)$   
 094 when  $\epsilon$  is sufficiently small, where  $N_{\mathcal{R}}(\epsilon)$  is the covering number of the reward function class and  
 095  $\eta$  is the KL-regularization coefficient.
- 096 • Then we showcase a novel analysis to upper bound the suboptimality gap of the KL-regularized  
 097 objective in contextual bandits, and propose a simple two-stage mixed sampling strategy for online  
 098 RLHF which achieves a sample complexity of  $\mathcal{O}(\max(\eta^2 D^2, \eta/\epsilon) \log N_{\mathcal{R}}(\epsilon/\delta))$  when the reward  
 099 scale is a constant, where  $D$  is the coverage coefficient of the reference policy  $\pi_0$  and  $\delta$  is the  
 100 confidence parameter. To the best of our knowledge, this is the first work to provide a sharp  
 101 sample complexity for KL-regularized contextual bandits.
- 102 • In Section 4, we extend our analysis to reinforcement learning from preference feedback. We  
 103 rigorously demonstrate that KL-regularization is essential for more efficient policy learning in  
 104 RLHF with preference data. We further propose a two-stage mixed sampling strategy for online  
 105 preference learning setting with a good coverage of  $\pi_0$ , which achieves a sample complexity of  
 106  $\mathcal{O}(\max(\eta^2 D^2, \eta/\epsilon) \log N_{\mathcal{R}}(\epsilon/\delta))$  when the reward scale is a constant.

## 2 PRELIMINARIES

In this section, we formally state the problem settings of reinforcement learning from human feedback (RLHF), where we consider two types of feedback: absolute rating and preference.

### 2.1 CONTEXTUAL BANDITS WITH KL REGULARIZATION

The first setting is the absolute-rating feedback, where we can query the ground-truth reward function to measure the quality of the responses by providing absolute reward value. For instance, in the NVIDIA Helpsteer project (Wang et al., 2023b; 2024c), human labelers are required to provide absolute score in five attributes: helpfulness, correctness, coherence, complexity, and verbosity. The dataset leads to many high-ranking open-source reward models, including the ArmoRM-Llama3-8B-v0.1 (Wang et al., 2024a;b), URM-LLaMa-3.1-8B<sup>1</sup>, and Llama-3.1-Nemotron-70B-Reward<sup>2</sup>. We also notice that recently this feedback framework is extended to other task such as video generation (He et al., 2024).

The absolute-rating feedback is directly modeled as reward functions (Wang et al., 2024a; Xiong et al., 2024b), and can be regarded as contextual bandits with KL regularization. In the contextual bandit setting, at each round  $t \geq 1$ , the agent observes a context  $x_t \in \mathcal{X}$  generated from a distribution  $d_0$  and chooses an action  $a_t \in \mathcal{A}$ . The agent receives a stochastic reward  $r_t \in \mathbb{R}$  that depends on the context  $x_t$  and the action  $a_t$ . The goal of the agent is to maximize the expected cumulative reward over  $T$  rounds.

The learner has access to a family of reward functions  $R(\theta, x, a)$  parameterized by  $\theta \in \Theta$ , such that there exists  $\theta_* \in \Theta$  satisfying  $\mathbb{E}[r_t | x_{1:t}, a_{1:t}] = R(\theta_*, x_t, a_t)$ . WLOG, we assume that the reward feedback  $r_t$  at all rounds is a non-negative real number bounded by  $B$ .

We consider a KL-regularized objective as follows:

$$Q(\pi) = \mathbb{E}_{x \sim d_0} \mathbb{E}_{a \sim \pi(\cdot|x)} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi(a|x)}{\pi_0(a|x)} \right], \quad (2.1)$$

where  $\pi_0$  is a known fixed policy, and  $\eta > 0$  is a hyperparameter that controls the trade-off between maximizing rewards and staying close to the reference policy  $\pi_0$ .

**Remark 2.1.** It is worth noting that entropy or Kullback-Leibler (KL) regularization is also widely used in contextual bandits (Berthet and Perchet, 2017; Wu et al., 2016) and deep reinforcement learning algorithms (Schulman et al., 2015; Fox et al., 2016; Schulman et al., 2017a; Haarnoja et al., 2017; 2018), where KL-divergence regularization is a popular technique for preventing drastic updates to the policy. Algorithms such as Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) explicitly incorporate KL-regularization to limit the policy updates during optimization, ensuring that the updated policy does not deviate too much from the current policy. This constraint promotes more stable and reliable learning, particularly in high-dimensional state-action spaces. Additionally, KL-regularization is central to Proximal Policy Optimization (PPO) (Schulman et al., 2017a), where a penalty term involving KL-divergence helps ensure updates remain within a ‘trust region’.

### 2.2 REINFORCEMENT LEARNING FROM PREFERENCE FEEDBACK

The second framework we consider is the preference feedback, which is a widely applied in projects such as Chat-GPT (OpenAI, 2023) and Claude (Bai et al., 2022). Specifically, when receiving a prompt  $x \in \mathcal{X}$ , and two actions (responses)  $a^1, a^2 \in \mathcal{A}$  from some LLM policy  $\pi(\cdot|x)$ , a preference oracle will give feedback  $y$  defined as follows:

**Definition 2.2** (Preference Oracle). A Preference Oracle is a function  $\mathbb{P} : \mathcal{X} \times \mathcal{A} \times \mathcal{A} \rightarrow \{0, 1\}$ . Given a context  $x \in \mathcal{X}$  and two actions  $a_1, a_2 \in \mathcal{A}$ , the oracle can be queried to obtain a preference signal  $y \sim \text{Bernoulli}(\mathbb{P}(x, a_1, a_2))$ , where  $y = 1$  indicates that  $a_1$  is preferred to  $a_2$  in the context  $x$ , and  $y = 0$  indicates the opposite.

To learn the preference, we follow Ouyang et al. (2022); Zhu et al. (2023); Rafailov et al. (2024); Liu et al. (2023); Xiong et al. (2024a) and assume that the preference oracle is measured by the

<sup>1</sup><https://huggingface.co/LxzGordon/URM-LLaMa-3.1-8B>

<sup>2</sup><https://huggingface.co/nvidia/Llama-3.1-Nemotron-70B-Reward>

162 difference of ground-truth reward functions  $R(\theta_*, x, a)$ , which is named the Bradley-Terry model  
 163 (Bradley and Terry, 1952b).

164 **Definition 2.3** (Bradley-Terry Model). The Bradley-Terry model is a probabilistic model for pair-  
 165 wise comparison data. Given a context  $x \in \mathcal{X}$  and two actions  $a_1, a_2 \in \mathcal{A}$ , the probability of  $a_1$   
 166 being preferred to  $a_2$  is modeled as

$$167 \mathbb{P}(x, a_1, a_2) = \frac{\exp(R(\theta_*, x, a_1))}{\exp(R(\theta_*, x, a_1)) + \exp(R(\theta_*, x, a_2))} = \sigma(R(\theta_*, x, a_1) - R(\theta_*, x, a_2)), \quad (2.2)$$

168 where  $\sigma(\cdot)$  is the sigmoid function.

169 The RLHF training always follows the fine-tuning process, which yields a reference policy  $\pi_0$ .  
 170 When performing RLHF on specific tasks, to avoid overfitting, we impose KL-regularization to the  
 171 learned reward model when optimizing the policy. Hence, our objective function is also (2.1).

### 172 2.3 ADDITIONAL NOTATIONS AND DEFINITIONS

173 In this subsection, we introduce the definitions shared by both settings.

174 **Reward function class.** We consider a function class  $\mathcal{R} = \{R(\theta, \cdot, \cdot) | \theta \in \Theta\}$  and for the realiz-  
 175 ability, we assume that the ground truth reward function  $R(\theta_*, x, a)$  is in the function class  $\mathcal{R}$ . Then,  
 176 we define the covering number of  $\mathcal{R}$  as follows.

177 **Definition 2.4** ( $\epsilon$ -cover and covering number). Given a function class  $\mathcal{F}$ , for each  $\epsilon > 0$ , an  $\epsilon$ -  
 178 cover of  $\mathcal{F}$  with respect to  $\|\cdot\|_\infty$ , denoted by  $\mathcal{C}(\mathcal{F}, \epsilon)$ , satisfies that for any  $f \in \mathcal{F}$ , we can find  
 179  $f' \in \mathcal{C}(\mathcal{F}, \epsilon)$  such that  $\|f - f'\|_\infty \leq \epsilon$ . The  $\epsilon$ -covering number, denoted as  $N_{\mathcal{F}}(\epsilon)$ , is the smallest  
 180 cardinality of such  $\mathcal{C}(\mathcal{F}, \epsilon)$ .

181 **Planning oracle.** Given a reward model, we can learn the policy by optimizing the KL-regularized  
 182 objective in (2.1). To simplify the analysis, we assume that there exists a planning oracle, which in  
 183 empirical can be efficiently approximated by rejection sampling (Liu et al., 2023), Gibbs sampling  
 184 (Xiong et al., 2024a), and iterative preference learning with a known reward (Dong et al., 2024).

185 **Definition 2.5** (Policy Improvement Oracle). For a reward function  $R(\theta, \cdot, \cdot) \in \mathcal{R}$  and a reference  
 186 policy  $\pi_0$ , for any prompt  $x \sim d_0$ , we can compute:

$$187 \pi_\theta^\eta(\cdot|x) := \operatorname{argmax}_{\pi(\cdot|x) \in \Delta(\mathcal{A})} \mathbb{E}_{a \sim \pi(\cdot|x)} \left[ R(\theta, x, a) - \frac{1}{\eta} \log \frac{\pi(a|x)}{\pi_0(a|x)} \right] \propto \pi_0(\cdot|x) \cdot \exp(\eta R(\theta, x, \cdot)).$$

188 Hence, the comparator policy is the solution of the oracle given the true reward function  $R(\theta^*, \cdot, \cdot)$ :  
 189  $\pi^*(\cdot|x) \propto \pi_0(\cdot|x) \cdot \exp(\eta R(\theta^*, \cdot, \cdot))$ . The **goal** is to minimize the sub-optimality of our learned policy  
 190  $\hat{\pi}$  with  $\pi^*: Q(\pi^*) - Q(\hat{\pi})$ .

191 **Coverage conditions.** It is crucial to assume that our data-collector policy  $\pi_0$  possesses good  
 192 coverage, which can ensure that the learned reward function can generalize well to unseen contexts  
 193 (prompts) and actions (responses), and thus can enable us to approximate the optimal policy. The  
 194 global coverage is the uniform cover over all the policies in the considered class  $\Pi$ , which is standard  
 195 in offline RL (Munos and Szepesvári, 2008; Song et al., 2024) and online RL (Xie et al., 2022; Rosset  
 196 et al., 2024). Essentially, Song et al. (2024) demonstrated that global coverage is necessary for  
 197 offline framework and Direct Preference Optimization (DPO) fails without global coverage. Hence,  
 198 we introduce two types of global coverage conditions.

199 **Definition 2.6** (Data Coverage). Given a reference policy  $\pi_0$ ,  $D^2$  is the minimum positive real  
 200 number satisfying  $\forall (x, a) \in \mathcal{X} \times \mathcal{A}$ , s.t.  $\pi(a|x) > 0$  and  $\forall b: \mathcal{X} \rightarrow [-B, B]$ , we have

$$201 \sup_{\theta, \theta' \in \Theta} \frac{|R(\theta', x, a) - R(\theta, x, a) - b(x)|^2}{\mathbb{E}_{x' \sim d_0} \mathbb{E}_{a' \sim \pi_0(\cdot|x')} |R(\theta', x', a') - R(\theta, x', a') - b(x')|^2} \leq D^2.$$

202 The coverage coefficient  $D$  measures how well the in-sample error induced by distribution  $d_0 \times$   
 203  $\pi_0$  can cover the out-of-sample error, identically speaking, it depicts the ability of  $\pi_0$  to cover  
 204 the action space. This concept is adapted from the F-design for online RL under general function  
 205 approximation (Agarwal et al.), and resembles the coverage coefficient for offline RL (Ye et al.,  
 206 2024c;a), and the eluder dimension (Wang et al., 2020; Ye et al., 2023; Agarwal et al., 2023) for  
 207 online RL.

**Definition 2.7** (Global-Policy Coverage). Given a reference policy  $\pi_0$ ,  $C_{\text{GL}}$  is the minimum positive real number satisfying that for any  $\pi : \mathcal{X} \rightarrow \mathcal{A}$

$$\sup_{x \sim d_0, a \in \mathcal{A}} \frac{\pi(a|x)}{\pi_0(a|x)} \leq C_{\text{GL}}.$$

The two conditions both require the reference policy to cover all possible policy distributions, which is standard and common in RL literature. Additionally, although the two conditions defined above are both global, it is obvious that  $D^2 \leq C_{\text{GL}}$ , indicating that it is more general to assume a finite  $D$  coefficient.

Because of the KL-regularization for RLHF, the learned policy will not move too far from the reference policy. Hence, it is natural to relax the global coverage to local coverage inside the KL-ball (Song et al., 2024).

**Definition 2.8** (Local KL-ball Coverage). Given a reference policy  $\pi_0$ , for a positive constant  $\rho_{\text{KL}} < \infty$ , and all policy satisfying that  $\mathbb{E}_{x \sim d_0}[\text{KL}(\pi, \pi_0)] \leq \rho_{\text{KL}}$ , we define

$$\sup_{x \sim d_0, a \in \mathcal{A}} \frac{\pi(a|x)}{\pi_0(a|x)} := C_{\rho_{\text{KL}}}.$$

**Remark 2.9** (Relation between Local and Global Coverage Conditions). The local coverage condition (Definition 2.8) is more precise because compared to the global conditions targeting all possible policies, it only constraint the coverage to a KL-ball. In Song et al. (2024), because of the specific form of the oracle (Definition 2.5), the considered policy class is  $\Pi = \{\pi(\cdot|\cdot) \propto \pi_0(\cdot|\cdot) \exp(\eta R(\theta, \cdot, \cdot)) : R(\theta, \cdot, \cdot) \in \mathcal{R}\}$ . Thus, they only need to assume that the condition hold for  $\rho = 2\eta B$ , indicating that  $C_{\rho_{\text{KL}}} \leq C_{\text{GL}}$ . On the other hand, the data coverage condition (Definition 2.6) is measured on the level of reward functions instead of policies. In this sense, the data coverage condition and local coverage condition do not encompass each other.

### 3 KL-REGULARIZED CONTEXTUAL BANDITS

#### 3.1 LOWER BOUND

In this section, we provide a lower bound for the KL-regularized contextual bandit problem.

**Theorem 3.1.** For any  $\epsilon \in (0, 1)$ ,  $\eta > 0$ , and any algorithm  $A$ , there exists a KL-regularized contextual bandit problem with  $O(1)$  coverage coefficient and reward function class  $\mathcal{R}$  such that  $A$  requires at least  $\Omega\left(\min\left(\frac{\eta \log N_{\mathcal{R}}(\epsilon)}{\epsilon}, \frac{\log N_{\mathcal{R}}(\epsilon)}{\epsilon^2}\right)\right)$  rounds to achieve a suboptimality gap of  $\epsilon$ .

**Remark 3.2.** The lower bound in Theorem 3.1 indicates that the sample complexity of the KL-regularized contextual bandit problem is  $\Omega(\eta \log N_{\mathcal{R}}(\epsilon)/\epsilon)$  when  $\epsilon$  is sufficiently small. In our proof, the KL-regularization term shifts the local landscape of the objective function, which prevents us to directly apply the standard bandit analysis, and thus requires a novel analysis to derive the new lower bound. This  $\Omega(\eta \log N_{\mathcal{R}}(\epsilon)/\epsilon)$  lower bound suggests that the KL-regularized contextual bandit problem enjoys a lower sample complexity compared to the standard contextual bandit problem.

#### 3.2 THE PROPOSED ALGORITHM

We present the algorithmic framework in Algorithm 1 for the KL-regularized contextual bandit problem, which serves as a theoretical model for online RLHF with absolute-rating feedback. The algorithm consists of two states:

- In the first stage, we sample  $m$  contexts (prompts) and actions (answers) from the foundation model  $\pi_0$  and observe the corresponding rewards (absolute ratings). These ratings can be regarded as noisy observations of the underlying reward function  $R(\theta_*, x, a)$ . In line 6, we compute an estimate of the reward function  $\hat{\theta}_0$  using least squares regression based on the collected data. In line 7, we apply the planning oracle to obtain the policy  $\pi_{\hat{\theta}_0}^\eta$  which maximizes the following KL-regularized estimated objective in Definition 2.5 with reward function  $R(\theta, \cdot, \cdot) = R(\hat{\theta}_0, \cdot, \cdot)$ .
- In the second stage, we utilize the trained policy  $\pi_{\hat{\theta}_0}^\eta$  to sample  $n$  contexts (prompts) and actions (responses). With the intermediate policy  $\pi_{\hat{\theta}_0}^\eta$ , we can collect new data  $\{(x_i, a_i, r_i)\}_{i=1}^n$  which is



**Algorithm 1** Two-Stage mixed-policy sampling

- 
- 1: **Input:**  $\eta, \epsilon, \pi_0, \Theta$ .  
 ▷ Use policy  $\pi_0$  to achieve sufficient data coverage
- 2: **for**  $i = 1, \dots, m$  **do**  
 3:   Sample context  $x_i^0 \sim d_0$  and action  $a_i^0 \sim \pi_0(\cdot|x_i^0)$ .  
 4:   Observe reward  $r_i^0 = R(\theta_*, x_i^0, a_i^0) + \epsilon_i^0$ , where  $\epsilon_i^0$  is the random noise.  
 5: **end for**
- 6: Compute the least square estimate of the reward function based on  $D_0 = \{(x_i^0, a_i^0, r_i^0)\}_{i=1}^m$ :
- $$\hat{\theta}_0 \leftarrow \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^m (R(\theta, x_i^0, a_i^0) - r_i^0)^2.$$
- 7: Apply the planning oracle to compute  $\pi_{\hat{\theta}_0}^\eta(\cdot) \propto \pi_0(\cdot) \exp(\eta R(\hat{\theta}_0, \cdot, \cdot))$ .  
 ▷ Use policy  $\pi_{\hat{\theta}_0}^\eta$  to sample new responses
- 8: **for**  $i = 1, \dots, n$  **do**  
 9:   Sample context  $x_i \sim d_0$  and action  $a_i \sim \pi_{\hat{\theta}_0}^\eta(\cdot|x_i)$ .  
 10:   Observe reward  $r_i = R(\theta_*, x_i, a_i) + \epsilon_i$ , where  $\epsilon_i$  is the random noise.  
 11: **end for**
- 12: Compute the least square estimate of the reward function using  $\{(x_i, a_i, r_i)\}_{i=1}^n$  together with  $D_0$ :
- $$\hat{\theta} \leftarrow \operatorname{argmin}_{\theta \in \Theta} \sum_{i=1}^m (R(\theta, x_i^0, a_i^0) - r_i^0)^2 + \sum_{i=1}^n (R(\theta, x_i, a_i) - r_i)^2.$$
- 13: **Output**  $\pi_{\hat{\theta}}^\eta(\cdot) \propto \pi_0(\cdot) \exp(\eta R(\hat{\theta}, \cdot, \cdot))$ .
- 

more aligned with the data distribution induced by the optimal policy  $\pi_*$ . In line 12, the algorithm combines data from both stages  $\{(x_i, a_i, r_i)\}_{i=1}^n$  and  $\{(x_i^0, a_i^0, r_i^0)\}_{i=1}^m$  to compute a refined least squares estimate  $\hat{\theta}$  of the reward function, minimizing the sum of squared errors across both datasets. By aggregating the two datasets together, there is an overlap between the data to compute  $\hat{\theta}$  and  $\hat{\theta}_0$ , so that the output policy  $\pi_{\hat{\theta}}^\eta$  is well covered by the intermediate policy  $\pi_{\hat{\theta}_0}^\eta$ .

## 3.3 THEORETICAL GUARANTEES

**Loose Bound of Previous Analysis.** The previous analysis is loose since they basically follow the techniques of bandits and neglect the significance of KL-regularization. For simplicity, We use shorthand notation  $R(\theta, x, \pi) = \mathbb{E}_{a \sim \pi(\cdot|x)} R(\theta, x, a)$  and denote  $\text{KL}(\pi(\cdot|x) \|\pi'(\cdot|x))$  by  $\text{KL}(\pi \|\pi')$  when there is no confusion. Estimator  $\hat{\theta}$  is estimated on a dataset  $\{(x_i, a_i, r_i) : x_i \sim d_0, s_i \sim \pi_0\}_{i=1}^n$ :  $\pi_{\hat{\theta}}^\eta = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{x \sim d_0} [R(\hat{\theta}, x, \pi) - \eta^{-1} \text{KL}(\pi \|\pi_0)]$ . The sub-optimality is decomposed as:

$$\begin{aligned} Q(\pi^*) - Q(\pi_{\hat{\theta}}^\eta) &= \mathbb{E}_{x \sim d_0} [R(\theta^*, x, \pi^*) - R(\hat{\theta}, x, \pi^*)] + \mathbb{E}_{x \sim d_0} [R(\hat{\theta}, x, \pi_{\hat{\theta}}^\eta) - R(\theta^*, x, \pi_{\hat{\theta}}^\eta)] \\ &\quad + \mathbb{E}_{x \sim d_0} [R(\hat{\theta}, x, \pi^*) - \eta^{-1} \text{KL}(\pi^* \|\pi_0)] - \mathbb{E}_{x \sim d_0} [R(\hat{\theta}, x, \pi_{\hat{\theta}}^\eta) - \eta^{-1} \text{KL}(\pi_{\hat{\theta}}^\eta \|\pi_0)] \\ &\leq \mathbb{E}_{x \sim d_0} [R(\theta^*, x, \pi^*) - R(\hat{\theta}, x, \pi^*) + R(\hat{\theta}, x, \pi_{\hat{\theta}}^\eta) - R(\theta^*, x, \pi_{\hat{\theta}}^\eta)], \end{aligned}$$

where the inequality holds since  $\pi_{\hat{\theta}}^\eta$  is the maximum.

Then, the suboptimality can be further bounded by using the coverage condition (Definition 2.7) and concentration inequalities: for any  $\pi \in \Pi$ , if  $n = \Theta(1/\epsilon^2)$ ,

$$\mathbb{E}_{x \sim d_0} \mathbb{E}_{a \sim \pi(\cdot|x)} [R(\theta^*, x, a) - R(\hat{\theta}, x, a)] \leq C_{\text{GL}} \mathbb{E}_{x \sim d_0} \mathbb{E}_{a \sim \pi_0(\cdot|x)} [R(\theta^*, x, a) - R(\hat{\theta}, x, a)] \leq C_{\text{GL}} \epsilon.$$

**Power of KL-regularization** The crucial point of the sharper result is utilizing the strong convexity of the objective  $Q$  because of the KL-regularization. Specifically, we take the first-order Taylor expansion of sub-optimality with respect to  $\{\Delta(x, a) = R(\hat{\theta}, x, a) - R(\theta^*, x, a) : a \in \mathcal{A}\}$

$$\begin{aligned} Q(\pi^*) - Q(\pi_{\hat{\theta}}^\eta) &= \eta \mathbb{E}_{x \sim d_0} \left[ \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \Delta^2(x, a) - \sum_{a_1, a_2 \in \mathcal{A}} \pi_f^\eta(a_1|x) \pi_f^\eta(a_2|x) \Delta(x, a_1) \Delta(x, a_2) \right] \\ &\leq \eta \mathbb{E}_{x \sim d_0} \left[ \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \Delta^2(x, a) \right], \end{aligned}$$

where  $f(\cdot, \cdot) = \gamma R(\hat{\theta}, \cdot, \cdot) + (1 - \gamma)R(\theta_*, \cdot, \cdot)$  ( $\gamma \in (0, 1)$ ) the inequality uses the fact that second term on the right-hand side of the equality is  $(\sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \Delta(x, a))^2 \geq 0$ .

Now, under Algorithm 1, the coverage condition (Definition 2.6) and with concentration inequalities, if the datsize  $m = \Theta(\eta^2 D^2 B^2)$ , we can prove that for  $\|R(\hat{\theta}, \cdot, \cdot) - R(\theta_*, \cdot, \cdot)\|_\infty \leq \eta^{-1}$  and  $\|R(\hat{\theta}_0, \cdot, \cdot) - R(\theta_*, \cdot, \cdot)\|_\infty \leq \eta^{-1}$ , which implies the whole-policy coverage condition:  $\|\pi_f^\eta(\cdot|x)/\pi_{\hat{\theta}_0}^\eta(\cdot|x)\|_\infty \leq e^4$ . Therefore, by setting  $n = \Theta(\eta/\epsilon)$ , we obtain that  $\pi_\theta^\eta$  is  $O(\epsilon)$  optimal.

The conclusion is presented in the following theorem.

**Theorem 3.3.** Suppose that Assumption 2.6 holds. For any  $\delta \in (0, 1/5)$ ,  $\epsilon > 0$  and constant  $c_{m,n} > 0$ , if we set  $m = \Theta(\eta^2 D^2 \cdot B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta))$  and  $n = \Theta(\eta/\epsilon \cdot B^2 \log(N_{\mathcal{R}}(\epsilon_c)/\delta))$  and  $\epsilon_c = \min\{\frac{\epsilon}{2(1+c_{m,n}^{-1})B}, \frac{1}{8(1+c_{m,n})B\eta^2 D^2}\}$ , then with probability at least  $1 - 5\delta$  the output policy of Algorithm 1  $\pi_\theta^\eta$  is  $O(\epsilon)$  optimal.

**Remark 3.4.** Theorem 3.3 shows that the sample complexity of Algorithm 1 is  $\mathcal{O}(\eta/\epsilon \log N_{\mathcal{R}}(\epsilon/\delta))$  when the reward scale is a constant and  $\epsilon$  is sufficiently small. The result indicates that the proposed two-stage mixed sampling strategy can achieve a suboptimality gap of  $\epsilon$  with only an additive dependence on the coverage coefficient  $D^2$ .

### 3.4 DISCUSSION: RESULT FOR LOCAL COVERAGE

In this subsection, we consider a weaker assumption as described in Definition 2.8.

**Corollary 3.5.** Let  $C_{\rho_{\text{KL}}}$  be defined in Definition 2.8 where  $\rho_{\text{KL}} = 2\eta B$ . For any  $\delta \in (0, 1/6)$  and  $\epsilon > 0$ , if we set  $n = c_{m,n} m = \Theta(C_{\rho_{\text{KL}}} \eta/\epsilon \cdot B \log(N_{\mathcal{R}}(\epsilon_c)/\delta))$  (where  $c_{m,n} > 0$  is a constant,  $\epsilon_c = \epsilon/(2(1+c_{m,n}^{-1})B)$ ) then with probability at least  $1 - 6\delta$  the output policy of Algorithm 2  $\pi_\theta^\eta$  is  $O(\epsilon)$  optimal.

*Proof of Corollary 3.5.* The proof follows the same lines as Theorem 4.3 by replacing the global coverage condition with the local coverage condition. It still holds that

$$Q(\pi^*) - Q(\pi_\theta^\eta) \leq \eta \cdot \mathbb{E}_{\pi_f^\eta} [(R(\hat{\theta}_0, x, a) - R(\theta_*, x, a))^2]$$

where  $\pi_f^\eta(a|x) \propto \pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))$  and  $f(\cdot, \cdot) = \gamma R(\hat{\theta}_0, \cdot, \cdot) + (1 - \gamma)R(\theta_*, \cdot, \cdot)$  for some  $\gamma \in (0, 1)$ . Thus, We have  $\text{KL}(\pi_f^\eta(\cdot|x) \|\pi_0(\cdot|x)) \leq 2\eta B$ , which further implies that

$$Q(\pi^*) - Q(\pi_\theta^\eta) \leq \eta \cdot C_{\rho_{\text{KL}}} \cdot O\left(\frac{1}{n} B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + B(1 + c_{m,n}^{-1})\epsilon_c\right)$$

by Lemma E.4. Then we can conclude by substituting the value of  $m$  into the suboptimality gap.  $\square$

## 4 REINFORCEMENT LEARNING FROM PREFERENCE FEEDBACK

In this section, we consider the problem of aligning the language model with preference feedback. As discussed in Section 2.2, at each round, we can sample a pair of actions (responses)  $a_1, a_2$  and call a preference oracle to get the preference label  $y \in \{0, 1\}$ , where  $y = 1$  means that the user prefers  $a_1$  over  $a_2$  (Definition 2.2).

Although preference feedback is believed to be more intuitive for human users and easier to collect, it also poses more challenges for the RLHF algorithms to effectively leverage the feedback signals since the reward signals are not directly observed.

In practice, RLHF with preference feedback typically involves

1. constructing a reward model based on the maximum likelihood estimation (MLE) of Bradley-Terry model from the preference feedback, and
2. applying RL algorithms like PPO (Schulman et al., 2017b) to train the language model so that it maximizes the reward signals with KL regularization (Ouyang et al., 2022; Bai et al., 2022; Touvron et al., 2023).

To analyze the above approach theoretically, we introduce the following assumption for step 1 to ensure the existence of an MLE estimation oracle which can globally maximize the likelihood of the Bradley-Terry model over all possible reward functions.

**Definition 4.1** (MLE estimation oracle). There exists an MLE estimation oracle that, given a set of context-action pairs  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$  generated from the Bradley-Terry model, can output the parameter  $\hat{\theta}$  such that

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \underbrace{y_i \cdot \log \sigma(R(\theta, x_i, a_i^1) - R(\theta, x_i, a_i^2)) + (1 - y_i) \cdot \log \sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1))}_{\mathcal{L}(\theta|x_i, a_i^1, a_i^2, y_i)}.$$

Following the previous analysis for RLHF (Xiong et al., 2024a), we assume the existence of a policy improvement oracle (Definition 2.5, corresponding to step 2) that can compute the optimal policy  $\pi_{\hat{\theta}}^{\eta}$  based on the reward function  $\hat{\theta}$ .

#### 4.1 LOWER BOUND

We provide a lower bound for the RLHF problem with preference feedback. The lower bound is derived by constructing a hard instance where the reward function is difficult to estimate from the preference feedback.

**Theorem 4.2.** For any  $\epsilon \in (0, 1)$ ,  $\eta > 0$ , and any algorithm  $A$ , there exists a KL-regularized preference learning problem as defined in Section 2.2 with  $O(1)$  coverage coefficient and reward function class  $\mathcal{R}$  such that  $A$  requires at least  $\Omega(\min(\frac{\eta \log N_{\mathcal{R}}(\epsilon)}{\epsilon}, \frac{\log N_{\mathcal{R}}(\epsilon)}{\epsilon^2}))$  samples to achieve a suboptimality gap of  $\epsilon$ .

#### 4.2 THEORETICAL GUARANTEES

We defer Algorithm 2, a 2-stage mixed-policy sampling algorithm for RLHF with preference feedback, to Appendix C for conciseness because of its similarity to Algorithm 1. We provide the theoretical guarantees for Algorithm 2 in the following theorem.

**Theorem 4.3.** Suppose that Assumption 2.6 holds. For any  $\delta \in (0, 1/6)$  and  $\epsilon > 0$ , if we set  $m = \Theta(\eta^2 D^2 \cdot e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta))$  and  $n = \Theta(\eta/\epsilon \cdot e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta))$  (where  $\epsilon_c = \min\{\frac{\epsilon}{2(1+c_{m,n}^{-1})e^B}, \frac{1}{8(1+c_{m,n})e^B \eta^2 D^2}\}$ ) then with probability at least  $1 - 6\delta$  the output policy of Algorithm 2  $\pi_{\hat{\theta}}^{\eta}$  is  $O(\epsilon)$  optimal.

**Remark 4.4** (Comparison with Hybrid Framework). We compare our two-stage mixed sampling method with hybrid frameworks. From the algorithmic perspective, a hybrid algorithm first learns from an offline dataset and then requires sufficient online iterations to ensure the performance (Xiong et al., 2024a). For example, for a finite action space with  $A$  actions, the number of online iterations should be  $\Theta(A)$ . In contrast, our method only requires two iterations of sampling from mixed policy and interacting with the environment. Moreover, the results of hybrid literature depend on both the coverage coefficient and the structure complexity of the function class (like the dimension for a linear function class or eluder dimension (Russo and Van Roy, 2013)). Our result only needs the coverage condition of the reference policy. More importantly, we obtain a sharper bound on the sample complexity and derive the additive dependence on the coverage coefficient.

**Remark 4.5.** Although the coefficient  $e^B$  appearing in sample size  $m, n$  can be exponentially large, this term is caused by the non-linearity of the link function for the preference model, and is common in RLHF literature (Zhu et al., 2023; Xiong et al., 2024a; Ye et al., 2024b; Song et al., 2024).

Theorem 4.3 shows that the sample complexity of Algorithm 2 is  $\mathcal{O}(\eta/\epsilon \log N_{\mathcal{R}}(\epsilon/\delta))$  when the reward scale is a constant and  $\epsilon$  is sufficiently small. The result indicates that the proposed two-stage mixed sampling strategy can achieve a suboptimality gap of  $\epsilon$  with only an additive dependence on the coverage coefficient  $D^2$ .

Besides, the algorithm only requires sampling from the reference policy  $\pi_0$  and the intermediate policy  $\pi_{\hat{\theta}_0}^{\eta}$ , which is more aligned with the practical scenarios where the preference feedback is collected from the human users and it is expensive to collect the data while the language model is being updated. Our result implies that we may achieve a near-optimal sample complexity by simply leveraging an intermediate policy to collect more data, and the training process of the reward model and the policy (language model) can be highly decoupled.



## REFERENCES

- 432  
433  
434 ABBASI-YADKORI, Y., PÁL, D. and SZEPESVÁRI, C. (2011). Improved algorithms for linear  
435 stochastic bandits. *Advances in neural information processing systems* **24**.  
436
- 437 ACHIAM, J., ADLER, S., AGARWAL, S., AHMAD, L., AKKAYA, I., ALEMAN, F. L., ALMEIDA,  
438 D., ALTENSCHMIDT, J., ALTMAN, S., ANADKAT, S. ET AL. (2023). Gpt-4 technical report.  
439 *arXiv preprint arXiv:2303.08774* .
- 440 AGARWAL, A., JIN, Y. and ZHANG, T. (2023). Vo  $q$  l: Towards optimal regret in model-free rl with  
441 nonlinear function approximation. In *The Thirty Sixth Annual Conference on Learning Theory*.  
442 PMLR.
- 443 AGARWAL, A., KAKADE, S. M., LEE, J. D. and MAHAJAN, G. (2020). Optimality and approx-  
444 imation with policy gradient methods in markov decision processes. In *Conference on Learning*  
445 *Theory*. PMLR.  
446
- 447 AGARWAL, A., KAKADE, S. M., LEE, J. D. and MAHAJAN, G. (2021). On the theory of policy  
448 gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learn-*  
449 *ing Research* **22** 1–76.
- 450 AGARWAL, A., QIAN, J., RAKHLIN, A. and ZHANG, T. (????). The non-linear  $f$ -design and  
451 applications to interactive learning. In *Forty-first International Conference on Machine Learning*.  
452
- 453 AHMED, Z., LE ROUX, N., NOROUZI, M. and SCHUURMANS, D. (2019). Understanding the  
454 impact of entropy on policy optimization. In *International conference on machine learning*.  
455 PMLR.
- 456 ANTHROPIC, A. (2023). Introducing claude.  
457
- 458 AZAR, M. G., GUO, Z. D., PIOT, B., MUNOS, R., ROWLAND, M., VALKO, M. and CALAN-  
459 DRIELLO, D. (2024). A general theoretical paradigm to understand learning from human prefer-  
460 ences. In *International Conference on Artificial Intelligence and Statistics*. PMLR.
- 461 BAI, J., BAI, S., CHU, Y., CUI, Z., DANG, K., DENG, X., FAN, Y., GE, W., HAN, Y., HUANG,  
462 F. ET AL. (2023). Qwen technical report. *arXiv preprint arXiv:2309.16609* .  
463
- 464 BAI, Y., JONES, A., NDOUSSE, K., ASKELL, A., CHEN, A., DASSARMA, N., DRAIN, D., FORT,  
465 S., GANGULI, D., HENIGHAN, T. ET AL. (2022). Training a helpful and harmless assistant with  
466 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862* .  
467
- 468 BERTHET, Q. and PERCHET, V. (2017). Fast rates for bandit optimization with upper-confidence  
469 frank-wolfe. *Advances in Neural Information Processing Systems* **30**.
- 470 BRADLEY, R. A. and TERRY, M. E. (1952a). Rank analysis of incomplete block designs: I. the  
471 method of paired comparisons. *Biometrika* **39** 324.  
472
- 473 BRADLEY, R. A. and TERRY, M. E. (1952b). Rank analysis of incomplete block designs: I. the  
474 method of paired comparisons. *Biometrika* **39** 324–345.
- 475 BROCKMAN, G. (2016). Openai gym. *arXiv preprint arXiv:1606.01540* .  
476
- 477 CALANDRIELLO, D., GUO, D., MUNOS, R., ROWLAND, M., TANG, Y., PIRES, B. A.,  
478 RICHEMOND, P. H., LAN, C. L., VALKO, M., LIU, T. ET AL. (2024). Human alignment of  
479 large language models through online preference optimisation. *arXiv preprint arXiv:2403.08635*  
480 .
- 481 CHEN, Z., DENG, Y., YUAN, H., JI, K. and GU, Q. (2024). Self-play fine-tuning converts weak  
482 language models to strong language models. *arXiv preprint arXiv:2401.01335* .  
483
- 484 CHRISTIANO, P. F., LEIKE, J., BROWN, T., MARTIC, M., LEGG, S. and AMODEI, D. (2017).  
485 Deep reinforcement learning from human preferences. *Advances in neural information process-*  
*ing systems* **30**.

- 486 COBBE, K., KOSARAJU, V., BAVARIAN, M., CHEN, M., JUN, H., KAISER, L., PLAPPERT, M.,  
487 TWOREK, J., HILTON, J., NAKANO, R. ET AL. (2021). Training verifiers to solve math word  
488 problems. *arXiv preprint arXiv:2110.14168* .
- 489  
490 DONG, H., XIONG, W., GOYAL, D., ZHANG, Y., CHOW, W., PAN, R., DIAO, S., ZHANG, J.,  
491 SHUM, K. and ZHANG, T. (2023). Raft: Reward ranked finetuning for generative foundation  
492 model alignment. *arXiv preprint arXiv:2304.06767* .
- 493  
494 DONG, H., XIONG, W., PANG, B., WANG, H., ZHAO, H., ZHOU, Y., JIANG, N., SAHOO, D.,  
495 XIONG, C. and ZHANG, T. (2024). Rlh workflow: From reward modeling to online rlhf. *arXiv  
496 preprint arXiv:2405.07863* .
- 497  
498 FOSTER, D. J., KAKADE, S. M., QIAN, J. and RAKHLIN, A. (2021). The statistical complexity of  
499 interactive decision making. *arXiv preprint arXiv:2112.13487* .
- 500  
501 FOX, R., PAKMAN, A. and TISHBY, N. (2016). Taming the noise in reinforcement learning via soft  
502 updates. In *32nd Conference on Uncertainty in Artificial Intelligence 2016, UAI 2016*. Association  
503 For Uncertainty in Artificial Intelligence (AUAI).
- 504  
505 GEIST, M., SCHERRER, B. and PIETQUIN, O. (2019). A theory of regularized markov decision  
506 processes. In *International Conference on Machine Learning*. PMLR.
- 507  
508 GOU, Z., SHAO, Z., GONG, Y., YELONG SHEN, YANG, Y., HUANG, M., DUAN, N. and CHEN,  
509 W. (2024). ToRA: A tool-integrated reasoning agent for mathematical problem solving. In *The  
510 Twelfth International Conference on Learning Representations*.  
511 URL <https://openreview.net/forum?id=Ep0TtjVoap>
- 512  
513 GUI, L., GÂRBACEA, C. and VEITCH, V. (2024). Bonbon alignment for large language models  
514 and the sweetness of best-of-n sampling. *arXiv preprint arXiv:2406.00832* .
- 515  
516 GULCEHRE, C., PAINE, T. L., SRINIVASAN, S., KONYUSHKOVA, K., WEERTS, L., SHARMA,  
517 A., SIDDHANT, A., AHERN, A., WANG, M., GU, C. ET AL. (2023). Reinforced self-training  
518 (rest) for language modeling. *arXiv preprint arXiv:2308.08998* .
- 519  
520 GUO, S., ZHANG, B., LIU, T., LIU, T., KHALMAN, M., LLINARES, F., RAME, A., MESNARD, T.,  
521 ZHAO, Y., PIOT, B. ET AL. (2024). Direct language model alignment from online ai feedback.  
522 *arXiv preprint arXiv:2402.04792* .
- 523  
524 HAARNOJA, T., TANG, H., ABBEEL, P. and LEVINE, S. (2017). Reinforcement learning with deep  
525 energy-based policies. In *International conference on machine learning*. PMLR.
- 526  
527 HAARNOJA, T., ZHOU, A., ABBEEL, P. and LEVINE, S. (2018). Soft actor-critic: Off-policy maximum  
528 entropy deep reinforcement learning with a stochastic actor. In *International conference  
529 on machine learning*. PMLR.
- 530  
531 HE, X., JIANG, D., ZHANG, G., KU, M., SONI, A., SIU, S., CHEN, H., CHANDRA, A., JIANG,  
532 Z., ARULRAJ, A. ET AL. (2024). Mantisscore: Building automatic metrics to simulate fine-  
533 grained human feedback for video generation. *arXiv preprint arXiv:2406.15252* .
- 534  
535 HENDRYCKS, D., BURNS, C., KADAVATH, S., ARORA, A., BASART, S., TANG, E., SONG, D.  
536 and STEINHARDT, J. (2021). Measuring mathematical problem solving with the math dataset.  
537 *arXiv preprint arXiv:2103.03874* .
- 538  
539 LIU, T., ZHAO, Y., JOSHI, R., KHALMAN, M., SALEH, M., LIU, P. J. and LIU, J. (2023). Sta-  
540 tistical rejection sampling improves preference optimization. *arXiv preprint arXiv:2309.06657*  
541 .
- 542  
543 MEI, J., XIAO, C., SZEPESVARI, C. and SCHUURMANS, D. (2020). On the global convergence  
544 rates of softmax policy gradient methods. In *International conference on machine learning*.  
545 PMLR.
- 546  
547 META, A. (2024). Introducing meta llama 3: The most capable openly available llm to date. *Meta  
548 AI* .

- 540 MUNOS, R. and SZEPESVÁRI, C. (2008). Finite-time bounds for fitted value iteration. *Journal of*  
541 *Machine Learning Research* **9**.
- 542
- 543 NAKANO, R., HILTON, J., BALAJI, S., WU, J., OUYANG, L., KIM, C., HESSE, C., JAIN, S.,  
544 KOSARAJU, V., SAUNDERS, W. ET AL. (2021). Webgpt: Browser-assisted question-answering  
545 with human feedback. *arXiv preprint arXiv:2112.09332* .
- 546
- 547 NEU, G., JONSSON, A. and GÓMEZ, V. (2017). A unified view of entropy-regularized markov  
548 decision processes. *arXiv preprint arXiv:1705.07798* .
- 549 OPENAI (2023). Gpt-4 technical report. *ArXiv* **abs/2303.08774**.
- 550
- 551 OUYANG, L., WU, J., JIANG, X., ALMEIDA, D., WAINWRIGHT, C., MISHKIN, P., ZHANG, C.,  
552 AGARWAL, S., SLAMA, K., RAY, A. ET AL. (2022). Training language models to follow in-  
553 structions with human feedback. *Advances in Neural Information Processing Systems* **35** 27730–  
554 27744.
- 555
- 556 RAFAILOV, R., SHARMA, A., MITCHELL, E., MANNING, C. D., ERMON, S. and FINN, C. (2024).  
557 Direct preference optimization: Your language model is secretly a reward model. *Advances in*  
558 *Neural Information Processing Systems* **36**.
- 559
- 560 ROSSET, C., CHENG, C.-A., MITRA, A., SANTACROCE, M., AWADALLAH, A. and XIE, T.  
561 (2024). Direct nash optimization: Teaching language models to self-improve with general pref-  
562 erences. *arXiv preprint arXiv:2404.03715* .
- 563
- 564 RUSSO, D. and VAN ROY, B. (2013). Eluder dimension and the sample complexity of optimistic  
565 exploration. *Advances in Neural Information Processing Systems* **26**.
- 566
- 567 SCHULMAN, J., CHEN, X. and ABBEEL, P. (2017a). Equivalence between policy gradients and  
568 soft q-learning. *arXiv preprint arXiv:1704.06440* .
- 569
- 570 SCHULMAN, J., LEVINE, S., ABBEEL, P., JORDAN, M. and MORITZ, P. (2015). Trust region  
571 policy optimization. In *International Conference on Machine Learning*. PMLR.
- 572
- 573 SCHULMAN, J., WOLSKI, F., DHARIWAL, P., RADFORD, A. and KLIMOV, O. (2017b). Proximal  
574 policy optimization algorithms. *arXiv preprint arXiv:1707.06347* .
- 575
- 576 SHANI, L., EFRONI, Y. and MANNOR, S. (2020). Adaptive trust region policy optimization: Global  
577 convergence and faster rates for regularized mdps. In *Proceedings of the AAAI Conference on*  
578 *Artificial Intelligence*, vol. 34.
- 579
- 580 SHAO, Z., WANG, P., ZHU, Q., XU, R., SONG, J., ZHANG, M., LI, Y., WU, Y. and GUO, D.  
581 (2024). Deepseekmath: Pushing the limits of mathematical reasoning in open language models.  
582 *arXiv preprint arXiv:2402.03300* .
- 583
- 584 SONG, Y., SWAMY, G., SINGH, A., BAGNELL, D. and SUN, W. (2024). The importance of online  
585 data: Understanding preference fine-tuning via coverage. In *ICML 2024 Workshop: Aligning*  
586 *Reinforcement Learning Experimentalists and Theorists*.
- 587
- 588 SUTTON, R. S. (2018). Reinforcement learning: An introduction. *A Bradford Book* .
- 589
- 590 SZEPESVÁRI, C. (2022). *Algorithms for reinforcement learning*. Springer nature.
- 591
- 592 TAJWAR, F., SINGH, A., SHARMA, A., RAFAILOV, R., SCHNEIDER, J., XIE, T., ERMON, S.,  
593 FINN, C. and KUMAR, A. (2024). Preference fine-tuning of llms should leverage suboptimal,  
594 on-policy data. *arXiv preprint arXiv:2404.14367* .
- 595
- 596 TEAM, G., ANIL, R., BORGEAUD, S., WU, Y., ALAYRAC, J.-B., YU, J., SORICUT, R., SCHALK-  
597 WYK, J., DAI, A. M., HAUTH, A. ET AL. (2023). Gemini: a family of highly capable multimodal  
598 models. *arXiv preprint arXiv:2312.11805* .
- 599
- 600 TONG, Y., ZHANG, X., WANG, R., WU, R. and HE, J. (2024). Dart-math: Difficulty-aware  
601 rejection tuning for mathematical problem-solving. *arXiv preprint arXiv:2407.13690* .

- 594 TOUVRON, H., MARTIN, L., STONE, K., ALBERT, P., ALMAHAIRI, A., BABAEI, Y., BASH-  
595 LYKOV, N., BATRA, S., BHARGAVA, P., BHOSALE, S. ET AL. (2023). Llama 2: Open foundation  
596 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* .  
597
- 598 WANG, H., LIN, Y., XIONG, W., YANG, R., DIAO, S., QIU, S., ZHAO, H. and ZHANG, T. (2024a).  
599 Arithmetic control of llms for diverse user preferences: Directional preference alignment with  
600 multi-objective rewards. *arXiv preprint arXiv:2402.18571* .  
601
- 602 WANG, H., XIONG, W., XIE, T., ZHAO, H. and ZHANG, T. (2024b). Interpretable preferences via  
603 multi-objective reward modeling and mixture-of-experts. *arXiv preprint arXiv:2406.12845* .  
604
- 605 WANG, R., SALAKHUTDINOV, R. R. and YANG, L. (2020). Reinforcement learning with gen-  
606 eral value function approximation: Provably efficient approach via bounded eluder dimension.  
607 *Advances in Neural Information Processing Systems* **33** 6123–6135.
- 608 WANG, Y., LIU, Q. and JIN, C. (2023a). Is rlhf more difficult than standard rl? a theoretical  
609 perspective. *Advances in Neural Information Processing Systems* **36** 76006–76032.
- 610 WANG, Z., DONG, Y., DELALLEAU, O., ZENG, J., SHEN, G., EGERT, D., ZHANG, J. J., SREED-  
611 HAR, M. N. and KUCHAIEV, O. (2024c). Helpsteer2: Open-source dataset for training top-  
612 performing reward models. *arXiv preprint arXiv:2406.08673* .  
613
- 614 WANG, Z., DONG, Y., ZENG, J., ADAMS, V., SREEDHAR, M. N., EGERT, D., DELALLEAU,  
615 O., SCOWCROFT, J. P., KANT, N., SWOPE, A. ET AL. (2023b). Helpsteer: Multi-attribute  
616 helpfulness dataset for steerlm. *arXiv preprint arXiv:2311.09528* .  
617
- 618 WU, J., OUYANG, L., ZIEGLER, D. M., STIENNON, N., LOWE, R., LEIKE, J. and CHRIS-  
619 TIANO, P. (2021). Recursively summarizing books with human feedback. *arXiv preprint*  
620 *arXiv:2109.10862* .
- 621 WU, Y., SHARIFF, R., LATTIMORE, T. and SZEPESVÁRI, C. (2016). Conservative bandits. In  
622 *International Conference on Machine Learning*. PMLR.  
623
- 624 XIE, T., FOSTER, D. J., BAI, Y., JIANG, N. and KAKADE, S. M. (2022). The role of coverage in  
625 online reinforcement learning. *arXiv preprint arXiv:2210.04157* .  
626
- 627 XIONG, W., DONG, H., YE, C., WANG, Z., ZHONG, H., JI, H., JIANG, N. and ZHANG, T.  
628 (2024a). Iterative preference learning from human feedback: Bridging theory and practice for  
629 rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*.
- 630 XIONG, W., SHI, C., SHEN, J., ROSENBERG, A., QIN, Z., CALANDRIELLO, D., KHALMAN, M.,  
631 JOSHI, R., PIOT, B., SALEH, M. ET AL. (2024b). Building math agents with multi-turn iterative  
632 preference learning. *arXiv preprint arXiv:2409.02392* .  
633
- 634 YE, C., HE, J., GU, Q. and ZHANG, T. (2024a). Towards robust model-based reinforcement  
635 learning against adversarial corruption. *arXiv preprint arXiv:2402.08991* .
- 636 YE, C., XIONG, W., GU, Q. and ZHANG, T. (2023). Corruption-robust algorithms with uncertainty  
637 weighting for nonlinear contextual bandits and markov decision processes. In *International Con-*  
638 *ference on Machine Learning*. PMLR.  
639
- 640 YE, C., XIONG, W., ZHANG, Y., JIANG, N. and ZHANG, T. (2024b). A theoretical analysis  
641 of nash learning from human feedback under general kl-regularized preference. *arXiv preprint*  
642 *arXiv:2402.07314* .
- 643 YE, C., YANG, R., GU, Q. and ZHANG, T. (2024c). Corruption-robust offline reinforcement learn-  
644 ing with general function approximation. *Advances in Neural Information Processing Systems*  
645 **36**.
- 646
- 647 YUE, Y., BRODER, J., KLEINBERG, R. and JOACHIMS, T. (2012). The k-armed dueling bandits  
problem. *Journal of Computer and System Sciences* **78** 1538–1556.

- ZHANG, X., XIONG, W., CHEN, L., ZHOU, T., HUANG, H. and ZHANG, T. (2024). From lists to emojis: How format bias affects model alignment. *arXiv preprint arXiv:2409.11704* .
- ZHAO, Y., JOSHI, R., LIU, T., KHALMAN, M., SALEH, M. and LIU, P. J. (2023). Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425* .
- ZHONG, H., FENG, G., XIONG, W., ZHAO, L., HE, D., BIAN, J. and WANG, L. (2024). Dpo meets ppo: Reinforced token optimization for rlhf. *arXiv preprint arXiv:2404.18922* .
- ZHU, B., JIAO, J. and JORDAN, M. I. (2023). Principled reinforcement learning with human feedback from pairwise or  $k$ -wise comparisons. *arXiv preprint arXiv:2301.11270* .
- ZIEBART, B. D., MAAS, A. L., BAGNELL, J. A., DEY, A. K. ET AL. (2008). Maximum entropy inverse reinforcement learning. In *Aaai*, vol. 8. Chicago, IL, USA.
- ZIEGLER, D. M., STIENNON, N., WU, J., BROWN, T. B., RADFORD, A., AMODEI, D., CHRISTIANO, P. and IRVING, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593* .

## A PREVIOUS UNDERSTANDING OF KL-REGULARIZATION IN RL

While we mainly focus on the theoretical understanding of KL-regularization in RLHF, it is also worth mentioning that our analysis for KL-regularized contextual bandits also contributes to the theoretical understanding the impact of KL-regularization in reinforcement learning since contextual bandits can be viewed as a simplified version of markov decision processes (MDPs).

In reinforcement learning, KL-regularization has been widely used to stabilize the learning process and prevent the policy from deviating too far from the reference policy. In this section, we provide a brief overview of the existing understanding of KL-regularization in decision-making problems. From the perspective of policy optimization, KL-regularization captures entropy regularization as a special case<sup>3</sup>, which is also an extensively used technique in reinforcement learning literature (Sutton, 2018; Szepesvári, 2022). There is a large body of literature that has explored the benefits of entropy regularization or KL-regularization in reinforcement learning (Schulman et al., 2015; Fox et al., 2016; Schulman et al., 2017a; Haarnoja et al., 2017; 2018; Ahmed et al., 2019). Most related to our work, Ahmed et al. (2019) provided a comprehensive understanding of the role of entropy regularization in reinforcement learning, showing that entropy regularization can improve the training efficiency and stability of the policy optimization process by changing the optimization landscape through experiments on continuous control tasks (Brockman, 2016).

Theoretically, Neu et al. (2017) provided a unified view of entropy regularization as approximate variants of Mirror Descent or Dual Averaging, and left the statistical justification for using entropy regularization in reinforcement learning as an open question. Geist et al. (2019) provided a framework for analyzing the error propagation in regularized MDPs, which also focused on the proof of the convergence for the policy optimization methods with regularization and lacked of a sharp sample complexity analysis.

## B OTHER RELATED LITERATURE

**Analyses for Policy Optimization with Regularization** While it is previously unknown whether regularization can improve the sample complexity of policy optimization without additional assumptions, there are some works that provided a sharp convergence rate in the presence of regularization (Mei et al., 2020; Shani et al., 2020; Agarwal et al., 2020; 2021). However, these works either assumed the access of exact or unbiased policy gradient or required uniform value function approximation error, which are not the standard case in sample-based reinforcement learning setting.

**RLHF Algorithms** There are mainly three types of RLHF algorithms: offline, online and hybrid. The most well-known offline algorithms are Slic (Zhao et al., 2023), Direct Preference Optimization (DPO) (Rafailov et al., 2024), Identity-PO (IPO) (Azar et al., 2024) and (SPIN) (Chen et al., 2024).

<sup>3</sup>We can regard the entropy regularization as a special case of KL-regularization by setting the reference policy as the uniform distribution.



They aim to approximate the closed-form solution of the optimization problem on a fixed offline dataset. For the online algorithms, the most representative one is Proximal Policy Optimization (PPO) (Schulman et al., 2017b). PPO has been used in the Chat-GPT (OpenAI, 2023), Gemini (Team et al., 2023), and Claude (Bai et al., 2022). However, the deep RL method PPO is known to be sample inefficient and unstable, making its success hard to reproduce for the open-source community. In response to this, there have been many efforts in proposing alternative algorithms to the PPO algorithm. The Reward ranked fine-tuning (RAFT) (also known as rejection sampling finetuning) (Dong et al., 2023; Touvron et al., 2023; Gulcehre et al., 2023; Gui et al., 2024) is a stable framework requiring minimal hyper-parameter tuning, which iteratively learns from the best-of-n policy (Nakano et al., 2021). This framework proves to be particularly effective in the reasoning task such as (Gou et al., 2024; Tong et al., 2024). However, the RAFT-like algorithms only use the positive signal by imitating the best-of-n sampling. To further improve the efficiency, there is an emerging body of literature that proposes online direct preference optimization by extending DPO or IPO to online iterative framework (Xiong et al., 2024a; Guo et al., 2024; Calandriello et al., 2024; Xiong et al., 2024b). Finally, for the third type, the common point of hybrid and online algorithms is that they both require further interaction with the preference oracle and on-policy data collection. The difference is that hybrid algorithms start with a pre-collected dataset (Xiong et al., 2024a; Song et al., 2024; Touvron et al., 2023), while the online algorithms learn from scratch.

**RLHF Theory** The theoretical study of RLHF can date back to the dueling bandit (Yue et al., 2012) and follow-up works on MDP (Wang et al., 2023a; Zhu et al., 2023). However, these works deviate from the practice because they do not realize the significance of KL-regularization and still choose the greedy policy that simply maximizes the reward. After this line of work, Xiong et al. (2024a); Ye et al. (2024b); Song et al. (2024) highlight the KL-regularization theoretically and incorporates the KL term into the learning objective. However, they circumvent the special advantages of KL-regularization and still follow the techniques in bandit analysis, thus obtaining a looser bound. In our paper, we establish a new lower bound and a sharper upper bound for the KL-regularized framework, thus validating the empirical advantage of KL-regularization. There are also some works extending KL-regularized RLHF from bandit problems to the Markov decision process (MDP) problems (Zhong et al., 2024; Xiong et al., 2024b). We expect that our techniques can also be extended to the MDP setting, which we leave for future work.

## C ALGORITHM FOR PREFERENCE FEEDBACK

In the first stage, we sample  $m$  context-action pairs  $\{(\tilde{x}_i, \tilde{a}_i^1, \tilde{a}_i^2, \tilde{y}_i)\}_{i=1}^m$  from the Bradley-Terry model and call the preference oracle to get the preference labels. We then compute the MLE estimator of the reward function  $\hat{\theta}_0$  based on the preference feedback in line 6. Afterwards, we apply the planning oracle to compute the optimal policy  $\pi_{\hat{\theta}_0}^*$  based on the reward function  $\hat{\theta}_0$  in line 7. Line 6 and line 7 correspond to the practical implementation of RLHF (Ouyang et al., 2022; Bai et al., 2022; Touvron et al., 2023) given a dataset of preference feedback.

In the second stage, we sample  $n$  context-action pairs  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$  using the intermediate policy  $\pi_{\hat{\theta}_0}^*$  and call the preference oracle to get the preference labels. We then compute the MLE estimator of the reward function  $\hat{\theta}$  based on the preference feedback from both stages. Finally, we apply the planning oracle to compute the optimal policy  $\pi_{\hat{\theta}}^*$  based on the reward function  $\hat{\theta}$ .

## D PROOFS FROM SECTION 3

### D.1 PROOF OF THEOREM 3.1

*Proof of Theorem 3.1.* Consider a simple case when  $|\mathcal{X}| = M$  and  $|\mathcal{A}| = 2$ . We suppose that the context  $x$  is drawn uniformly from  $\mathcal{X}$  at the beginning of each round. Let  $\Theta$  be the set consisting of mappings from  $\mathcal{X}$  to  $\mathcal{A} = \{0, 1\}$ . For each  $\theta \in \Theta$ , we have  $R(\theta, x, a) = \begin{cases} 1/2 + c & \text{if } a = \theta(x), \\ 1/2 & \text{if } a \neq \theta(x), \end{cases}$  where  $c > 0$  is a constant, and  $\theta(x)$  is the optimal action under context  $x$  when the model is  $\theta$ .

For any  $(\theta, x, a) \in \Theta \times \mathcal{X} \times \mathcal{A}$ , we assume the reward feedback  $r \sim \text{Bernoulli}(R(\theta, x, a))$  when the model is  $\theta$  and  $a$  is chosen under context  $x$ .

**Algorithm 2** Stage mixed-policy sampling for preference feedback

- 
- 1: **Input:**  $\eta, \epsilon, \pi_0, \Theta$ .  
 ▷ Use policy  $\pi_0$  to achieve sufficient data coverage
- 2: **for**  $i = 1, \dots, m$  **do**
- 3:     Sample context  $\tilde{x}_i \sim d_0$  and 2 actions  $\tilde{a}_i^1, \tilde{a}_i^2 \sim \pi_0(\cdot|\tilde{x}_i)$ .
- 4:     Observe preference label  $\tilde{y}_i \in \{0, 1\}$  from the preference oracle defined in Definition 2.2.
- 5: **end for**
- 6: Compute the MLE estimator of the reward function based on  $\{(\tilde{x}_i, \tilde{a}_i^1, \tilde{a}_i^2, \tilde{y}_i)\}_{i=1}^m$ :
- $$\hat{\theta}_0 \leftarrow \operatorname{argmax}_{\theta} \sum_{i=1}^m \tilde{y}_i \cdot \log \sigma(R(\theta, \tilde{x}_i, \tilde{a}_i^1) - R(\theta, \tilde{x}_i, \tilde{a}_i^2)) + (1 - \tilde{y}_i) \cdot \log \sigma(R(\theta, \tilde{x}_i, \tilde{a}_i^2) - R(\theta, \tilde{x}_i, \tilde{a}_i^1)).$$
- 7: Apply the planning oracle to compute  $\pi_{\hat{\theta}_0}^{\eta}(\cdot|\cdot) \propto \pi_0(\cdot|\cdot) \exp(\eta R(\hat{\theta}_0, \cdot, \cdot))$ .  
 ▷ Use policy  $\pi_{\hat{\theta}_0}^{\eta}$  to sample new responses
- 8: **for**  $i = 1, \dots, n$  **do**
- 9:     Sample context  $x_i \sim d_0$  and 2 actions  $a_i^1, a_i^2 \sim \pi_{\hat{\theta}_0}^{\eta}(\cdot|x_i)$ .
- 10:     Observe preference label  $y_i \in \{0, 1\}$  from the preference oracle defined in Definition 2.2.
- 11: **end for**
- 12: Compute the MLE estimator of the reward function using  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$  together with  $\{(\tilde{x}_i, \tilde{a}_i^1, \tilde{a}_i^2, \tilde{y}_i)\}_{i=1}^m$ :
- $$\hat{\theta} \leftarrow \operatorname{argmax}_{\theta} \sum_{i=1}^m \tilde{y}_i \cdot \log \sigma(R(\theta, \tilde{x}_i, \tilde{a}_i^1) - R(\theta, \tilde{x}_i, \tilde{a}_i^2)) + (1 - \tilde{y}_i) \cdot \log \sigma(R(\theta, \tilde{x}_i, \tilde{a}_i^2) - R(\theta, \tilde{x}_i, \tilde{a}_i^1))$$
- $$+ \sum_{i=1}^n y_i \cdot \log \sigma(R(\theta, x_i, a_i^1) - R(\theta, x_i, a_i^2)) + (1 - y_i) \cdot \log \sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1))$$
- 13: **Output**  $\pi_{\hat{\theta}}^{\eta}(\cdot|\cdot) \propto \pi_0(\cdot|\cdot) \exp(\eta R(\hat{\theta}, \cdot, \cdot))$ .
- 

We pick a pair of model  $\theta_1, \theta_2$  in  $\Theta$ , such that  $\theta_1(x) = \begin{cases} \theta_2(x) & \text{if } x \neq x_0, \\ 1 - \theta_2(x) & \text{if } x = x_0. \end{cases}$

We denote by  $\mathbb{P}_{\theta}, \mathbb{E}_{\theta}$  the probability measure and expectation under the model  $\theta$ .

Applying Pinsker's inequality (Lemma F.3), we have for all event  $A$  measurable with respect to the filtration generated by the observations,

$$|\mathbb{P}_{\theta_1}(A) - \mathbb{P}_{\theta_2}(A)| \leq \sqrt{\frac{1}{2} \log(1 - 4c^2) \mathbb{E}_{\theta_1}[N(x_0)]} \leq \sqrt{2c^2 \mathbb{E}_{\theta_1}[N(x_0)]} = \sqrt{2c^2 T/M},$$

where the first inequality follows from the chain rule of KL divergence, and the fact that  $\mathbb{E}_{\theta_1}[N(x_0)] = T/M$ .

Set  $A$  to be the event that  $\pi_{out}(\theta_1(x_0)|x_0) > 1/2$ . Then we have

$$\mathbb{P}_{\theta_1}(\pi_{out}(\theta_1(x_0)|x_0) \leq 1/2) + \mathbb{P}_{\theta_2}(\pi_{out}(\theta_2(x_0)|x_0) \leq 1/2) \geq 1 - |\mathbb{P}_{\theta_1}(A) - \mathbb{P}_{\theta_2}(A)| \geq 1 - \sqrt{2c^2 T/M}.$$

If the model  $\theta$  is uniformly drawn from  $\Theta$ , then we have

$$\mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{P}_{\theta}(\pi_{out}(\theta(x_0)) \leq 1/2) \geq \frac{1}{2} - \sqrt{c^2 T/2M}$$

for an arbitrary  $x_0$ .

Then we consider the following suboptimality gap:

$$\begin{aligned} & \mathbb{E}_{\pi_{\theta_*}^{\eta}} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\theta_*}^{\eta}(a|x)}{\pi_0(a|x)} \right] - \mathbb{E}_{\pi_{out}} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{out}(a|x)}{\pi_0(a|x)} \right] \\ &= \frac{1}{\eta} \mathbb{E}_{\pi_{\theta_*}^{\eta}} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\theta_*}^{\eta}(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{out}(a|x)} \right] \end{aligned}$$

810

811

$$= \frac{1}{\eta} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right],$$

812

where the last equality follows from the fact that  $\pi_{\theta_*}^\eta \propto \pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))$ .

813

To bound the suboptimality gap, we further have

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

$$\begin{aligned} & \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \\ &= \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \frac{1}{M} \sum_{x \in \mathcal{X}} \mathbb{E}_{a \sim \pi_{out}(\cdot|x)} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \\ &\geq \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \frac{1}{M} \sum_{x \in \mathcal{X}} \mathbb{P}_\theta(\pi_{out}(\theta(x)) \leq 1/2) \cdot \left[ \frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} \right] \\ &\geq \left( \frac{1}{2} - \sqrt{c^2 T / 2M} \right) \left[ \frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} \right] \end{aligned} \quad (\text{D.1})$$

Note that

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

$$\begin{aligned} & \frac{d}{du} \left[ \frac{1}{2} \ln \frac{1 + e^{-u}}{2} + \frac{1}{2} \ln \frac{1 + e^u}{2} \right] \Big|_{u=0} = \frac{1}{2} \left[ \frac{1}{1 + \exp(-u)} - \frac{1}{1 + \exp(u)} \right] \Big|_{u=0} = 0, \\ & \frac{d^2}{du^2} \left[ \frac{1}{2} \ln \frac{1 + e^{-u}}{2} + \frac{1}{2} \ln \frac{1 + e^u}{2} \right] = \frac{\exp(u)}{[1 + \exp(u)]^2}. \end{aligned}$$

Thus, applying Taylor's expansion on the right-hand side of (D.1), we have

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

$$\mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \geq \frac{1}{2} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 2M} \right) \eta^2 c^2 \cdot \frac{1}{3 + \exp(\eta c)}$$

When  $\epsilon < 1/64\eta$ , we can set  $c = 8\sqrt{\epsilon/\eta}$ . To achieve a suboptimality gap of  $\epsilon$ , we need to satisfy:

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

$$\frac{1}{2} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 2M} \right) \eta^2 c^2 \cdot \frac{1}{3 + \exp(\eta c)} \leq \eta \epsilon,$$

indicating that  $T \geq \frac{\eta M}{512\epsilon} = \Omega\left(\frac{\eta M}{\epsilon}\right)$ .

When  $\epsilon \geq 1/64\eta$ , or equivalently,  $\eta \geq 1/64\epsilon$ , we employ a different lower bound for (D.1) as follows:

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

$$\begin{aligned} \frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} &= \frac{1}{2} \ln \frac{2 + \exp(\eta c) + \exp(-\eta c)}{4} \\ &\geq \frac{1}{2} \cdot \frac{1}{2} \left( \ln \frac{\exp(\eta c) + \exp(-\eta c)}{2} \right) \\ &\geq \frac{1}{4} (\eta c - \ln 2), \end{aligned} \quad (\text{D.2})$$

where the first inequality follows from Jensen's inequality.

Substituting (D.2) into (D.1), we have

852

853

854

855

856

857

858

859

860

861

862

863

$$\epsilon \geq \frac{1}{\eta} \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \geq \frac{1}{4} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 2M} \right) (\eta c - \ln 2) \cdot \frac{1}{\eta}.$$

Set  $c = 64\epsilon$ . Then we have  $T = \Omega(M/\epsilon^2)$ . □

## D.2 PROOF OF THEOREM 3.3

We start with the following lemma, which provides an on-policy generalization bound for the reward function. Due to the on-policy nature of the algorithm (i.e., the usage of intermediate  $\pi_{\theta_0}^\eta$ ), we can leverage the covering number of the reward function class  $\mathcal{R}$  to derive the generalization error. Since we are using a fixed policy  $\pi_{\theta_0}^\eta$  to sample in the second stage, we can derive the generalization error of the reward function as follows:

**Lemma D.1** (Generalization error of reward function). For an arbitrary policy  $\pi$ , a set of context-action pairs  $\{(x_i, a_i)\}_{i=1}^n$  generated i.i.d. from  $\pi$ , and a distance threshold  $0 < \epsilon_c \leq B$ , we have with probability at least  $1 - \delta$ , for any pair of parameters  $\theta_1$  and  $\theta_2$ ,

$$\begin{aligned} & \mathbb{E}_\pi |R(\theta_1, x, a) - R(\theta_2, x, a)|^2 \\ & \leq \frac{2}{n} \sum_{i=1}^n |R(\theta_1, x_i, a_i) - R(\theta_2, x_i, a_i)|^2 + \frac{32B^2}{3n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 10\epsilon_c B. \end{aligned}$$

*Proof.* We first consider an  $\epsilon_c$ -net  $\mathcal{R}^c$  of the reward function class  $\mathcal{R}$  where  $\mathcal{R}^c = \{R(\theta, \cdot, \cdot) | \theta \in \Theta^c\}$  with size  $N_{\mathcal{R}}(\epsilon_c)$ . For any  $R(\theta, \cdot, \cdot) \in \mathcal{R}$ , there exists  $\theta^c$  such that  $\|R(\theta, \cdot, \cdot) - R(\theta^c, \cdot, \cdot)\|_\infty \leq \epsilon_c$ .

By Lemma F.1, for each pair of  $\theta_1^c, \theta_2^c \in \Theta^c$  (corresponding to  $\theta_1, \theta_2$ ), we have with probability at least  $1 - \delta$ ,

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n (R(\theta_1^c, x_i, a_i) - R(\theta_2^c, x_i, a_i))^2 - \mathbb{E}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2 \right| \\ & \leq \sqrt{\frac{2\text{Var}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2}{n} \log(2/\delta)} + \frac{2}{3n} B^2 \log(2/\delta) \\ & \leq \sqrt{\frac{2B^2 \mathbb{E}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2}{n} \log(2/\delta)} + \frac{2}{3n} B^2 \log(2/\delta) \end{aligned}$$

where the second inequality follows from the fact that  $R(\theta_1^c, x, a), R(\theta_2^c, x, a) \leq B$ .

Using union bound over all  $\theta_1^c, \theta_2^c \in \Theta^c$ , we have with probability at least  $1 - \delta$ , for all  $\theta_1^c, \theta_2^c \in \Theta^c$ ,

$$\begin{aligned} & \mathbb{E}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2 - \frac{1}{n} \sum_{i=1}^n (R(\theta_1^c, x_i, a_i) - R(\theta_2^c, x_i, a_i))^2 \\ & \leq \sqrt{\frac{4B^2 \mathbb{E}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2}{n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)} + \frac{4B^2}{3n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta), \end{aligned}$$

from which we further obtain the following inequality by Lemma F.2,

$$\mathbb{E}_\pi |R(\theta_1^c, x, a) - R(\theta_2^c, x, a)|^2 \leq \frac{2}{n} \sum_{i=1}^n (R(\theta_1^c, x_i, a_i) - R(\theta_2^c, x_i, a_i))^2 + \frac{32B^2}{3n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta). \quad (\text{D.3})$$

Then we can complete the proof by the definition of  $\epsilon$ -net.  $\square$

Next, we provide the following lemma, which gives an upper bound on the cumulative square error of the learned reward function.

**Lemma D.2** (Confidence bound for reward function). For an arbitrary policy  $\pi$ , and a set of data  $\{(x_i, a_i, r_i)\}_{i=1}^n$  generated i.i.d. from  $\pi$ , suppose that  $\hat{\theta}$  is the least squares estimator of  $\theta_*$ , i.e.,  $\hat{\theta} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n (R(\theta, x_i, a_i) - r_i)^2$ . Then for any threshold  $\epsilon_c > 0$ , with probability at least  $1 - \delta$ , it holds that

$$\sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))^2 \leq 16B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 4\epsilon_c n B.$$

*Proof.* We have the following inequality for  $\sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))^2$ ,

$$\begin{aligned} & \sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))^2 \\ & = \sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - r_i)^2 - \sum_{i=1}^n (R(\theta_*, x_i, a_i) - r_i)^2 \end{aligned}$$

$$\begin{aligned}
& + 2 \sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))(r_i - R(\theta_*, x_i, a_i)) \\
& \leq 2 \sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))(r_i - R(\theta_*, x_i, a_i)),
\end{aligned}$$

where the last inequality follows from the fact that  $\sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - r_i)^2 \leq \sum_{i=1}^n (R(\theta_*, x_i, a_i) - r_i)^2$ .

We then consider an  $\epsilon_c$ -net  $\mathcal{R}^c$  of the reward function class  $\mathcal{R}$  where  $\mathcal{R}^c = \{R(\theta, \cdot, \cdot) | \theta \in \Theta^c\}$  with size  $N_{\mathcal{R}}(\epsilon_c)$ . For any  $R(\theta, \cdot, \cdot) \in \mathcal{R}$ , there exists  $\theta^c$  such that  $\|R(\theta, x, a) - R(\theta^c, x, a)\|_{\infty} \leq \epsilon_c$ . From Azuma-Hoeffding inequality, with probability at least  $1 - \delta$ , it holds for all  $\theta \in \Theta^c$  that

$$\begin{aligned}
& \sum_{i=1}^n (R(\theta, x_i, a_i) - R(\theta_*, x_i, a_i))(r_i - R(\theta_*, x_i, a_i)) \\
& \leq \sqrt{2B^2 \sum_{i=1}^n (R(\theta, x_i, a_i) - R(\theta_*, x_i, a_i))^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)}.
\end{aligned}$$

Then we further have with probability at least  $1 - \delta$ , there exists  $\|R(\theta^c, \cdot, \cdot) - R(\hat{\theta}, \cdot, \cdot)\| \leq \epsilon_c$  such that

$$\begin{aligned}
& \sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))(r_i - R(\theta_*, x_i, a_i)) \\
& \leq \sqrt{2B^2 \sum_{i=1}^n (R(\theta, x_i, a_i) - R(\theta_*, x_i, a_i))^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 2\epsilon_c n B},
\end{aligned}$$

which implies that

$$\sum_{i=1}^n (R(\hat{\theta}, x_i, a_i) - R(\theta_*, x_i, a_i))^2 \leq 16B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 4\epsilon_c n B \quad (\text{D.4})$$

from Lemma F.2.  $\square$

With the above lemmas, we are now ready to prove the following lemma that bounds the estimation error of the reward function  $R(\hat{\theta}, \cdot, \cdot)$  under the sampled policy  $\pi_{\hat{\theta}_0}^{\eta}$ .

**Lemma D.3.** Let  $\hat{\theta}_0$  be the least squares estimator of the reward function based on the data  $\{(x_i^0, a_i^0, r_i^0)\}_{i=1}^m$  generated from  $\pi_0$  as defined in Algorithm 1. Then for any threshold  $\epsilon_c > 0$ , with probability at least  $1 - 2\delta$ , we have

$$\mathbb{E}_{\pi_{\hat{\theta}_0}^{\eta}} |R(\hat{\theta}, x, a) - R(\theta_*, x, a)|^2 \leq \frac{43B^2}{n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 10\epsilon_c(1 + m/n)B.$$

*Proof.* By Lemma D.1, we have with probability at least  $1 - \delta$ , the following upper bound holds for  $\mathbb{E}_{\pi_{\hat{\theta}_0}^{\eta}} |R(\theta_1, x, a) - R(\theta_2, x, a)|^2$ ,

$$\begin{aligned}
& \mathbb{E}_{\pi_{\hat{\theta}_0}^{\eta}} |R(\theta_1, x, a) - R(\theta_2, x, a)|^2 \\
& \leq \frac{2}{n} \sum_{i=1}^n |R(\theta_1, x_i, a_i) - R(\theta_2, x_i, a_i)|^2 + \frac{32B^2}{3n} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 10\epsilon_c B. \quad (\text{D.5})
\end{aligned}$$

By Lemma D.2, with probability at least  $1 - \delta$

$$\sum_{i=1}^n |R(\theta_*, x_i, a_i) - R(\hat{\theta}, x_i, a_i)|^2 \leq 16B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 4\epsilon_c(n + m)B. \quad (\text{D.6})$$

Then we can complete the proof using a union bound and substituting (D.6) into (D.5).  $\square$



**Lemma D.4.** If  $m \geq 128\eta^2 D^2 B^2 \cdot \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)$ , and there exists a positive constant  $c_{m,n} > 0$  such that  $n = c_{m,n}m$  in Algorithm 1 and Assumption 2.6 holds, then by taking  $\epsilon_c \leq \min\{B, (8(1 + c_{m,n})B\eta^2 D^2)^{-1}\}$ , with probability at least  $1 - 3\delta$ , we have

$$\eta|R(\hat{\theta}_0, x, a) - R(\theta_*, x, a)| \leq 1, \quad \eta|R(\hat{\theta}, x, a) - R(\theta_*, x, a)| \leq 1$$

for any pair  $(x, a) \in \mathcal{X} \times \mathcal{A}$  such that  $\pi_0(a|x) > 0$ .

*Proof.* By Lemma D.1, with probability at least  $1 - \delta$ , for all  $\theta_1, \theta_2 \in \Theta$ , we have

$$\mathbb{E}_{\pi_0} |R(\theta_1, x, a) - R(\theta_2, x, a)|^2 \leq \frac{2}{m} \sum_{i=1}^m |R(\theta_1, x_i^0, a_i^0) - R(\theta_2, x_i^0, a_i^0)|^2 + \frac{32B^2}{3m} \log(2N_{\mathcal{R}}(\epsilon_c)/\delta).$$

By Lemma D.2, with probability at least  $1 - \delta$ , we have

$$\sum_{i=1}^m |R(\hat{\theta}_0, x_i^0, a_i^0) - R(\theta_*, x_i^0, a_i^0)|^2 \leq 16B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 4\epsilon_c m.$$

Also, with probability at least  $1 - \delta$ , we have

$$\sum_{i=1}^m |R(\theta_*, x_i^0, a_i^0) - R(\hat{\theta}, x_i^0, a_i^0)|^2 \leq 16B^2 \log(2N_{\mathcal{R}}(\epsilon_c)/\delta) + 4\epsilon_c(m+n)B.$$

Similar to the proof of Lemma D.3, we have if  $m \geq 128\eta^2 D^2 B^2 \cdot \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)$ ,  $n = c_{m,n}m$ , then with probability at least  $1 - 3\delta$ ,

$$\mathbb{E}_{\pi_0} |R(\theta_*, x, a) - R(\hat{\theta}_0, x, a)|^2 \leq 1/\eta^2 D^2, \quad \mathbb{E}_{\pi_0} |R(\theta_*, x, a) - R(\hat{\theta}, x, a)|^2 \leq 1/\eta^2 D^2.$$

which implies that  $\eta|R(\hat{\theta}_0, x, a) - R(\theta_*, x, a)| \leq 1$  and  $\eta|R(\hat{\theta}, x, a) - R(\theta_*, x, a)| \leq 1$  for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$  such that  $\pi_0(a|x) > 0$ .  $\square$

*Proof of Theorem 3.3.* We have

$$\begin{aligned} & \mathbb{E}_{\pi_{\theta_*}^\eta} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\theta_*}^\eta(a|x)}{\pi_0(a|x)} \right] - \mathbb{E}_{\pi_{\hat{\theta}}^\eta} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\hat{\theta}}^\eta(a|x)}{\pi_0(a|x)} \right] \\ &= \frac{1}{\eta} \mathbb{E}_{\pi_{\theta_*}^\eta} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\theta_*}^\eta(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{\hat{\theta}}^\eta} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\hat{\theta}}^\eta(a|x)} \right] \\ &= \frac{1}{\eta} \mathbb{E}_{x \sim d_0} [\ln Z_{\theta_*}^\eta(x)] - \frac{1}{\eta} \mathbb{E}_{x \sim d_0} [\ln Z_{\hat{\theta}}^\eta(x)] - \mathbb{E}_{x \sim d_0} \left[ \sum_{a \in \mathcal{A}} \pi_{\hat{\theta}}^\eta(a|x) \cdot (R(\theta_*, x, a) - R(\hat{\theta}, x, a)) \right] \end{aligned}$$

For an arbitrary reward function  $f : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ , let  $\Delta(x, a) = f(x, a) - R(\theta_*, x, a)$ . Consider the following first derivative of  $J(f) = \ln Z_f^\eta(x) - \eta \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \Delta(x, a)$ , where  $Z_f^\eta(x) = \sum_{a \in \mathcal{A}} \pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))$  and  $\pi_f^\eta(a|x) \propto \pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))$ .

$$\begin{aligned} & \frac{\partial}{\partial \Delta(x, a)} \left[ \ln Z_f^\eta(x) - \eta \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \Delta(x, a) \right] \\ &= \frac{1}{Z_f^\eta(x)} \cdot \pi_0(a|x) \exp(\eta \cdot f(x, a)) \cdot \eta - \eta \cdot \pi_f^\eta(a|x) \\ & \quad - \eta \cdot \Delta(x, a) \cdot \frac{\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))}{Z_f^\eta(x)} \cdot \eta + \eta \cdot \Delta(x, a) \cdot \frac{[\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))]^2}{[Z_f^\eta(x)]^2} \cdot \eta \\ & \quad + \eta \sum_{a' \in \mathcal{A} \setminus \{a\}} \frac{\pi_0(a'|x) \cdot \exp(\eta \cdot f(x, a'))}{Z_f^\eta(x)} \cdot \eta \cdot \Delta(x, a') \cdot \frac{\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))}{Z_f^\eta(x)} \end{aligned}$$

$$= -\eta^2 \pi_f^\eta(a|x) \Delta(x, a) + \eta^2 [\pi_f^\eta(a|x)]^2 \cdot \Delta(x, a) + \eta^2 \sum_{a' \in \mathcal{A} \setminus \{a\}} \pi_f^\eta(a'|x) \pi_f^\eta(a|x) \Delta(x, a').$$

Therefore, there exists  $f(\cdot, \cdot) = \gamma R(\hat{\theta}, \cdot, \cdot) + (1 - \gamma) R(\theta_*, \cdot, \cdot)$  such that  $(\gamma \in (0, 1))$

$$\begin{aligned} \mathbb{E}_{x \sim d_0} [J(R(\hat{\theta}, \cdot, \cdot)) - J(R(\theta_*, \cdot, \cdot))] &= \frac{1}{\eta} \mathbb{E}_{x \sim d_0} \left[ -\eta^2 \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \gamma \cdot (R(\hat{\theta}, x, a) - R(\theta_*, x, a))^2 \right] \\ &+ \frac{1}{\eta} \mathbb{E}_{x \sim d_0} \left[ \gamma \eta^2 \sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \pi_f^\eta(a_1|x) \pi_f^\eta(a_2|x) (R(\hat{\theta}, x, a_1) - R(\theta_*, x, a_1)) (R(\hat{\theta}, x, a_2) - R(\theta_*, x, a_2)) \right] \\ &\geq -\eta \cdot \mathbb{E}_{\pi_f^\eta} [(R(\hat{\theta}, x, a) - R(\theta_*, x, a))^2] \end{aligned}$$

From Lemma D.4, if  $m \geq 128\eta^2 D^2 B^2 \cdot \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)$ , for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$  such that  $\pi_0(a|x) > 0$ , it holds that

$$\eta |R(\hat{\theta}_0, x, a) - R(\theta_*, x, a)| \leq 1, \quad \eta |R(\hat{\theta}, x, a) - R(\theta_*, x, a)| \leq 1,$$

which means that for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$

$$\frac{\pi_f^\eta(a|x)}{\pi_{\hat{\theta}_0}^\eta(a|x)} \leq e^4.$$

Let  $\epsilon_c = \min\{\frac{\epsilon}{(1+c_{m,n})B}, \frac{1}{8(1+c_{m,n})B\eta^2 D^2}, B\}$ . From Lemma D.3, if  $m \geq 128\eta^2 D^2 B^2 \cdot \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)$  and  $n \geq \eta/\epsilon \cdot B^2 \log(N_{\mathcal{R}}(\epsilon_c)/\delta)$  and  $n = c_{m,n}m$  then with high probability the output policy  $\pi_\theta^\eta$  is  $O(\epsilon)$  optimal.  $\square$

## E PROOFS FROM SECTION 4

### E.1 PROOF OF THEOREM 4.2

*Proof of Theorem 4.2.* The proof follows a similar construction as the one for Theorem 3.1. Consider a simple case when  $|\mathcal{X}| = M$  and  $|\mathcal{A}| = 2$ . We suppose that the context  $x$  is drawn uniformly from  $\mathcal{X}$  at the beginning of each round. Let  $\Theta$  be the set consisting of mappings from  $\mathcal{X}$  to  $\mathcal{A} = \{0, 1\}$ . For each  $\theta \in \Theta$ , we have  $R(\theta, x, a) = \begin{cases} c & \text{if } a = \theta(x), \\ 0 & \text{if } a \neq \theta(x), \end{cases}$  where  $c > 0$  is a constant, and  $\theta(x)$  is the optimal action under context  $x$  when the model is  $\theta$ .

We pick a pair of model  $\theta_1, \theta_2$  in  $\Theta$ , such that  $\theta_1(x) = \begin{cases} \theta_2(x) & \text{if } x \neq x_0, \\ 1 - \theta_2(x) & \text{if } x = x_0. \end{cases}$

We denote by  $\mathbb{P}_\theta, \mathbb{E}_\theta$  the probability measure and expectation under the model  $\theta$ .

Applying Pinsker's inequality (Lemma F.3), we have for all event  $A$  measurable with respect to the filtration generated by the observations,

$$|\mathbb{P}_{\theta_1}(A) - \mathbb{P}_{\theta_2}(A)| \leq \sqrt{\log(1/2 + e^c/4 + e^{-c}/4) \mathbb{E}_{\theta_1}[N(x_0)]} \leq \sqrt{c^2 \mathbb{E}_{\theta_1}[N(x_0)]} = \sqrt{c^2 T/M},$$

where the first inequality follows from the chain rule of KL divergence, and the fact that  $\mathbb{E}_{\theta_1}[N(x_0)] = T/M$ .

Set  $A$  to be the event that  $\pi_{out}(\theta_1(x_0)|x_0) > 1/2$ . Then we have

$$\mathbb{P}_{\theta_1}(\pi_{out}(\theta_1(x_0)|x_0) \leq 1/2) + \mathbb{P}_{\theta_2}(\pi_{out}(\theta_2(x_0)|x_0) \leq 1/2) \geq 1 - |\mathbb{P}_{\theta_1}(A) - \mathbb{P}_{\theta_2}(A)| \geq 1 - \sqrt{c^2 T/M}.$$

If the model  $\theta$  is uniformly drawn from  $\Theta$ , then we have

$$\mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{P}_\theta(\pi_{out}(\theta(x_0)) \leq 1/2) \geq \frac{1}{2} - \sqrt{c^2 T/4M}$$

for an arbitrary  $x_0$ .

Then we consider the following suboptimality gap:

$$\begin{aligned}
& \mathbb{E}_{\pi_{\theta_*}^\eta} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\theta_*}^\eta(a|x)}{\pi_0(a|x)} \right] - \mathbb{E}_{\pi_{out}} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{out}(a|x)}{\pi_0(a|x)} \right] \\
&= \frac{1}{\eta} \mathbb{E}_{\pi_{\theta_*}^\eta} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\theta_*}^\eta(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{out}(a|x)} \right] \\
&= \frac{1}{\eta} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right],
\end{aligned}$$

where the last equality follows from the fact that  $\pi_{\theta_*}^\eta \propto \pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))$ .

To bound the suboptimality gap, we further have

$$\begin{aligned}
& \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \\
&= \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \frac{1}{M} \sum_{x \in \mathcal{X}} \mathbb{E}_{a \sim \pi_{out}(\cdot|x)} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \\
&\geq \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \frac{1}{M} \sum_{x \in \mathcal{X}} \mathbb{P}_\theta(\pi_{out}(\theta(x)) \leq 1/2) \cdot \left[ \frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} \right] \\
&\geq \left( \frac{1}{2} - \sqrt{c^2 T / 4M} \right) \left[ \frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} \right] \tag{E.1}
\end{aligned}$$

Note that

$$\begin{aligned}
& \frac{d}{du} \left[ \frac{1}{2} \ln \frac{1 + e^{-u}}{2} + \frac{1}{2} \ln \frac{1 + e^u}{2} \right] \Big|_{u=0} = \frac{1}{2} \left[ \frac{1}{1 + \exp(-u)} - \frac{1}{1 + \exp(u)} \right] \Big|_{u=0} = 0, \\
& \frac{d^2}{du^2} \left[ \frac{1}{2} \ln \frac{1 + e^{-u}}{2} + \frac{1}{2} \ln \frac{1 + e^u}{2} \right] = \frac{\exp(u)}{[1 + \exp(u)]^2}.
\end{aligned}$$

Thus, applying Taylor's expansion on the right-hand side of (E.1), we have

$$\mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \geq \frac{1}{2} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 4M} \right) \eta^2 c^2 \cdot \frac{1}{3 + \exp(\eta c)}$$

When  $\epsilon < 1/64\eta$ , we can set  $c = 8\sqrt{\epsilon/\eta}$ . To achieve a suboptimality gap of  $\epsilon$ , we need to satisfy:

$$\frac{1}{2} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 4M} \right) \eta^2 c^2 \cdot \frac{1}{3 + \exp(\eta c)} \leq \eta \epsilon,$$

indicating that  $T \geq \frac{\eta M}{512\epsilon} = \Omega\left(\frac{\eta M}{\epsilon}\right)$ .

When  $\epsilon \geq 1/64\eta$ , or equivalently,  $\eta \geq 1/64\epsilon$ , we employ a different lower bound for (D.1) as follows:

$$\begin{aligned}
\frac{1}{2} \ln \frac{1 + \exp(-\eta c)}{2} + \frac{1}{2} \ln \frac{1 + \exp(\eta c)}{2} &= \frac{1}{2} \ln \frac{2 + \exp(\eta c) + \exp(-\eta c)}{4} \\
&\geq \frac{1}{2} \cdot \frac{1}{2} \left( \ln \frac{\exp(\eta c) + \exp(-\eta c)}{2} \right) \\
&\geq \frac{1}{4} (\eta c - \ln 2), \tag{E.2}
\end{aligned}$$

where the first inequality follows from Jensen's inequality.

Substituting (E.2) into (E.1), we have

$$\epsilon \geq \frac{1}{\eta} \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\pi_{out}} \left[ \ln \frac{\pi_{out}(a|x)}{\pi^*(a|x)} \right] \geq \frac{1}{4} \cdot \left( \frac{1}{2} - \sqrt{c^2 T / 4M} \right) (\eta c - \ln 2) \cdot \frac{1}{\eta}.$$

Set  $c = 64\epsilon$ . Then we have  $T = \Omega(M/\epsilon^2)$ .

□

## E.2 PROOF OF THEOREM 4.3

First, we provide the following lemma for the connection between the likelihood loss and the reward difference, which is a key step to upper bound the reward difference between  $\hat{\theta}$  and  $\theta_*$ .

**Lemma E.1.** For an arbitrary policy  $\pi$ , and a set of context-action pairs  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$  generated i.i.d. from the Bradley-Terry model and  $\pi$ , we have with probability at least  $1 - \delta$ , for any  $s \leq n$ ,

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^s \mathcal{L}(\theta | x_i, a_i^1, a_i^2, y_i) - \mathcal{L}(\theta_* | x_i, a_i^1, a_i^2, y_i) \\ & \leq \log(1/\delta) - \frac{1}{8} e^{-B} \sum_{i=1}^s ([R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 \end{aligned}$$

*Proof.* Applying Lemma F.4 to the sequence

$$\left\{ -\frac{1}{2} y_i \cdot \log \frac{\sigma(R(\theta_*, x_i, a_i^1) - R(\theta_*, x_i, a_i^2))}{\sigma(R(\theta, x_i, a_i^1) - R(\theta, x_i, a_i^2))} - \frac{1}{2} (1-y_i) \cdot \log \frac{\sigma(R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1))}{\sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1))} \right\}_{i=1}^n,$$

We have with probability at least  $1 - \delta$ , for all  $s \leq n$ ,

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^s \mathcal{L}(\theta | x_i, a_i^1, a_i^2, y_i) - \mathcal{L}(\theta_* | x_i, a_i^1, a_i^2, y_i) \\ & \leq \log(1/\delta) + \sum_{i=1}^s \log \left( \sqrt{\sigma(R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)) \cdot \sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1))} \right. \\ & \quad \left. + \sqrt{\sigma(R(\theta_*, x_i, a_i^1) - R(\theta_*, x_i, a_i^2)) \cdot \sigma(R(\theta, x_i, a_i^1) - R(\theta, x_i, a_i^2))} \right) \\ & = \log(1/\delta) - \frac{1}{2} \sum_{i=1}^s \left( \sqrt{\sigma(R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1))} - \sqrt{\sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1))} \right)^2 \\ & \leq \log(1/\delta) - \frac{1}{8} \sum_{i=1}^s (\sigma(R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)) - \sigma(R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1)))^2 \\ & \leq \log(1/\delta) - \frac{1}{8} e^{-B} \sum_{i=1}^s ([R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2, \end{aligned}$$

where the equality follows from the fact that  $\sigma(r) + \sigma(-r) = 1$  and the last inequality holds since  $\sigma'(r) = \sigma(r) \cdot (1 - \sigma(r)) \geq e^{-B}$  for all  $r \in [-B, B]$ .  $\square$

To further control the error bound for the reward function with the help of Lemma E.1, we introduce the following lemma to show that the likelihood function class  $\mathcal{L}$  can be well-covered by the  $\epsilon$ -net of the reward function class  $\mathcal{R}$ .

**Lemma E.2** (Covering number of  $\mathcal{L}$ ). For any  $\epsilon_c > 0$ , consider an  $\epsilon_c$ -net  $\mathcal{R}^c = \{R(\theta, \cdot, \cdot) | \theta \in \Theta^c\}$  for the reward function class  $\mathcal{R}$  with size  $N_{\mathcal{R}}(\epsilon_c)$ . Then for any  $\theta \in \Theta$ , there exists  $\theta^c \in \Theta^c$  such that for any  $s \in [n]$ ,

$$\sum_{i=1}^s \mathcal{L}(\theta | x_i, a_i^1, a_i^2, y_i) \leq \sum_{i=1}^s \mathcal{L}(\theta^c | x_i, a_i^1, a_i^2, y_i) + 2s\epsilon_c.$$

*Proof.* For any  $r \in \mathbb{R}$ , we have

$$\frac{d \log(\sigma(r))}{dr} = \frac{1}{\sigma(r)} \cdot \sigma(r) \cdot (1 - \sigma(r)) = 1 - \sigma(r) \in (0, 1).$$

Thus, for any  $\theta \in \Theta$ ,  $x \in \mathcal{X}$ ,  $a^1, a^2 \in \mathcal{A}$  and  $y \in \{0, 1\}$ , there exists  $\theta^c \in \Theta^c$  such that

$$|\mathcal{L}(\theta | x, a^1, a^2, y) - \mathcal{L}(\theta^c | x, a^1, a^2, y)|$$

$$\leq |[R(\theta, x, a^1) - R(\theta, x, a^2)] - [R(\theta^c, x, a^1) - R(\theta^c, x, a^2)]| = 2\epsilon_c.$$

1188  
1189  
1190  
1191

□

1192 With the above two lemmas, we are now ready to provide the confidence bound for the MLE esti-  
1193 mator of the reward function.

1194 **Lemma E.3.** Consider a set of context-action pairs  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$  where labels  $\{y_i\}_{i=1}^n$  are  
1195 generated independently from the Bradley-Terry model. Suppose that  $\hat{\theta}$  is the MLE estimator as  
1196 defined in Definition 4.1. We have with probability at least  $1 - \delta$ ,

$$\sum_{i=1}^n ([R(\hat{\theta}, x_i, a_i^2) - R(\hat{\theta}, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 \leq O(e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B n\epsilon_c).$$

1197  
1198  
1199

1200 *Proof.* By Lemma E.1 and Lemma E.2, we have with probability at least  $1 - \delta$ , for any  $\theta \in \Theta$ ,

1201

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^n \mathcal{L}(\theta|x_i, a_i^1, a_i^2, y_i) - \mathcal{L}(\theta_*|x_i, a_i^1, a_i^2, y_i) \\ & \leq \log(N_{\mathcal{R}}(\epsilon_c)/\delta) - \frac{1}{8} e^{-B} \sum_{i=1}^n ([R(\theta, x_i, a_i^2) - R(\theta, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 + O(n\epsilon_c). \end{aligned}$$

1202  
1203  
1204

1205 Since  $\hat{\theta}$  is the MLE estimator, we have  $\sum_{i=1}^n \mathcal{L}(\hat{\theta}|x_i, a_i^1, a_i^2, y_i) - \mathcal{L}(\theta_*|x_i, a_i^1, a_i^2, y_i) \geq 0$ , which  
1206 further implies

1207  
1208  
1209

$$0 \leq \log(N_{\mathcal{R}}(\epsilon_c)/\delta) - \frac{1}{8} e^{-B} \sum_{i=1}^n ([R(\hat{\theta}, x_i, a_i^2) - R(\hat{\theta}, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 + O(n\epsilon_c).$$

1210  
1211  
1212

1213 Then we have

1214

$$\sum_{i=1}^n ([R(\hat{\theta}, x_i, a_i^2) - R(\hat{\theta}, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 \leq O(e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B n\epsilon_c).$$

1215  
1216  
1217

□

1218 Finally, we provide the on-policy confidence bound for the squared reward difference between the  
1219 MLE estimator  $\hat{\theta}$  and the optimal reward function  $\theta_*$ .

1220  
1221

1222 **Lemma E.4.** Consider an arbitrary policy  $\pi$ , and a set of context-action pairs  $\{(x_i, a_i^1, a_i^2, y_i)\}_{i=1}^n$   
1223 generated i.i.d. from the Bradley-Terry model and  $\pi$ . Suppose that  $\hat{\theta}$  is the MLE estimator. We have  
1224 with probability at least  $1 - 2\delta$ , there exists a mapping  $b : \mathcal{X} \rightarrow \mathbb{R}$  such that

1225  
1226  
1227

$$\mathbb{E}_{\pi} [(R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b(x))^2] \leq O\left(\frac{1}{n} e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B \epsilon_c\right).$$

1228  
1229

1230 *Proof.* By Lemma E.3, we have with probability at least  $1 - \delta$ ,

1231

$$\sum_{i=1}^n ([R(\hat{\theta}, x_i, a_i^2) - R(\hat{\theta}, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2 \leq O(e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B n\epsilon_c).$$

1232  
1233

1234 We consider an  $\epsilon_c$ -net  $\mathcal{R}^c = \{R(\theta, \cdot, \cdot) | \theta \in \Theta^c\}$  for the reward function class  $\mathcal{R}$  with size  $N_{\mathcal{R}}(\epsilon_c)$ .  
1235 For any  $R(\theta, \cdot, \cdot)$ , there exists  $R(\theta^c, \cdot, \cdot)$  such that

1236  
1237

$$|R(\theta, x, a) - R(\theta^c, x, a)| \leq O(\epsilon_c)$$

1238

1239 for all  $x \in \mathcal{X}, a \in \mathcal{A}$ .

1240

1241 Applying Lemma F.1, with probability at least  $1 - \delta$ , we have

1242

$$\sum_{i=1}^n ([R(\theta^c, x_i, a_i^2) - R(\theta^c, x_i, a_i^1)] - [R(\theta_*, x_i, a_i^2) - R(\theta_*, x_i, a_i^1)])^2$$

1243  
1244



$$\begin{aligned}
& -n\mathbb{E}_{x\sim d_0}\mathbb{E}_{a^1,a^2\sim\pi}[(R(\theta^c,x,a^1)-R(\theta_*,x,a^1)-R(\theta^c,x,a^2)+R(\theta_*,x,a^2))^2] \\
& \leq \sqrt{\sum_{i=1}^n 4B^2\mathbb{E}_{x\sim d_0}\mathbb{E}_{a^1,a^2\sim\pi}[(R(\theta^c,x,a^1)-R(\theta_*,x,a^1)-R(\theta^c,x,a^2)+R(\theta_*,x,a^2))^2] \log(N_{\mathcal{R}}(\epsilon_c)/\delta)} \\
& \quad + \frac{8}{3}B^2\log(N_{\mathcal{R}}(\epsilon_c)/\delta)
\end{aligned}$$

for all  $\theta^c \in \Theta^c$ .

From Lemma F.2 and the definition of  $\Theta^c$ , we further have

$$\begin{aligned}
& \mathbb{E}_{x\sim d_0}\mathbb{E}_{a^1,a^2\sim\pi}[(R(\widehat{\theta},x,a^1)-R(\theta_*,x,a^1)-R(\widehat{\theta},x,a^2)+R(\theta_*,x,a^2))^2] \tag{E.3} \\
& \leq O\left(\frac{1}{n}B^2\log(N_{\mathcal{R}}(\epsilon_c)/\delta) + \frac{1}{n}\sum_{i=1}^n ([R(\widehat{\theta},x_i,a_i^2)-R(\widehat{\theta},x_i,a_i^1)] - [R(\theta_*,x_i,a_i^2)-R(\theta_*,x_i,a_i^1)])^2 + B\epsilon_c\right), \tag{E.4}
\end{aligned}$$

from which we can further derive that

$$\begin{aligned}
& \mathbb{E}_{x\sim d_0}\mathbb{E}_{a^1,a^2\sim\pi}[(R(\widehat{\theta},x,a^1)-R(\theta_*,x,a^1)-R(\widehat{\theta},x,a^2)+R(\theta_*,x,a^2))^2] \\
& \leq O\left(\frac{1}{n}e^B\log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B\epsilon_c\right)
\end{aligned}$$

with probability at least  $1 - 2\delta$  from Lemma E.3 and the union bound.

We can then complete the proof by setting

$$b(x) = \mathbb{E}_{a^2\sim\pi(\cdot|x)}[R(\widehat{\theta},x,a^2) - R(\theta_*,x,a^2)].$$

□

**Lemma E.5** (Coverage of  $\pi_*$  and  $\pi_{\widehat{\theta}}$  by  $\pi_{\widehat{\theta}_0}$ ). If  $m \geq 32\eta^2 D^2 e^B \log(N_{\mathcal{R}}(\epsilon_c))$ ,  $n = c_{m,n}m$  and  $\epsilon_c \leq \frac{1}{(1+c_{m,n})e^B\eta^2 D^2}$  in Algorithm 2 and Assumption 2.6 holds, then with probability at least  $1 - 4\delta$ , there exists  $b_1 : \mathcal{X} \rightarrow \mathbb{R}$  and  $b_2 : \mathcal{X} \rightarrow \mathbb{R}$  such that

$$\eta|R(\widehat{\theta}_0,x,a) - R(\theta_*,x,a) - b_1(x)| \leq 1, \quad \eta|R(\widehat{\theta},x,a) - R(\theta_*,x,a) - b_2(x)| \leq 1$$

for all  $x \in \mathcal{X}$ ,  $a \in \mathcal{A}$  such that  $\pi_0(a|x) > 0$ .

*Proof.* By Lemma E.3 and the union bound, we have with probability at least  $1 - \delta$ , it holds that

$$\begin{aligned}
& \sum_{i=1}^m ([R(\widehat{\theta},\tilde{x}_i,\tilde{a}_i^2) - R(\widehat{\theta},\tilde{x}_i,\tilde{a}_i^1)] - [R(\theta_*,\tilde{x}_i,\tilde{a}_i^2) - R(\theta_*,\tilde{x}_i,\tilde{a}_i^1)])^2 \\
& \quad + \sum_{i=1}^n ([R(\widehat{\theta},x_i,a_i^2) - R(\widehat{\theta},x_i,a_i^1)] - [R(\theta_*,x_i,a_i^2) - R(\theta_*,x_i,a_i^1)])^2 \\
& \leq O(e^B\log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B(n+m)\epsilon_c). \tag{E.5}
\end{aligned}$$

Consider an  $\epsilon_c$ -net  $\mathcal{R}^c = \{R(\theta, \cdot, \cdot) | \theta \in \Theta^c\}$  for the reward function class  $\mathcal{R}$  with size  $N_{\mathcal{R}}(\epsilon_c)$ . For any  $R(\theta, \cdot, \cdot)$ , there exists  $R(\theta^c, \cdot, \cdot)$  such that

$$|R(\theta, x, a) - R(\theta^c, x, a)| \leq O(\epsilon_c)$$

for all  $x \in \mathcal{X}$ ,  $a \in \mathcal{A}$ .

Applying Lemma F.1, with probability at least  $1 - \delta$ , we have

$$\begin{aligned}
& \sum_{i=1}^m ([R(\theta^c,\tilde{x}_i,\tilde{a}_i^2) - R(\theta^c,\tilde{x}_i,\tilde{a}_i^1)] - [R(\theta_*,x_i,a_i^2) - R(\theta_*,x_i,a_i^1)])^2 \\
& \quad - m\mathbb{E}_{x\sim d_0}\mathbb{E}_{a^1,a^2\sim\pi_0}[(R(\theta^c,x,a^1) - R(\theta_*,x,a^1) - R(\theta^c,x,a^2) + R(\theta_*,x,a^2))^2]
\end{aligned}$$

$$\begin{aligned} &\leq \sqrt{\sum_{i=1}^m 4B^2 \mathbb{E}_{x \sim d_0} \mathbb{E}_{a^1, a^2 \sim \pi_0} [(R(\theta^c, x, a^1) - R(\theta_*, x, a^1) - R(\theta^c, x, a^2) + R(\theta_*, x, a^2))^2]} \log(N_{\mathcal{R}}(\epsilon_c)/\delta) \\ &\quad + \frac{8}{3} B^2 \log(N_{\mathcal{R}}(\epsilon_c)/\delta) \end{aligned}$$

for all  $\theta^c \in \Theta^c$ .

From Lemma F.2 and the definition of  $\Theta^c$ , we further have

$$\begin{aligned} &\mathbb{E}_{x \sim d_0} \mathbb{E}_{a^1, a^2 \sim \pi} [(R(\hat{\theta}, x, a^1) - R(\theta_*, x, a^1) - R(\hat{\theta}, x, a^2) + R(\theta_*, x, a^2))^2] \\ &\leq O\left(\frac{1}{m} B^2 \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + \frac{1}{m} \sum_{i=1}^n ([R(\hat{\theta}, \tilde{x}_i, \tilde{a}_i^2) - R(\hat{\theta}, \tilde{x}_i, \tilde{a}_i^1)] - [R(\theta_*, \tilde{x}_i, \tilde{a}_i^2) - R(\theta_*, \tilde{x}_i, \tilde{a}_i^1)])^2 + B\epsilon_c\right). \end{aligned} \tag{E.6}$$

Substituting (E.5) into (E.6), we have with probability at least  $1 - 2\delta$ ,

$$\begin{aligned} &\mathbb{E}_{x \sim d_0} \mathbb{E}_{a^1, a^2 \sim \pi_0} [(R(\hat{\theta}, x, a^1) - R(\theta_*, x, a^1) - R(\hat{\theta}, x, a^2) + R(\theta_*, x, a^2))^2] \\ &\leq O\left(\frac{1}{m} e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B \cdot \frac{n+m}{m} \cdot \epsilon_c\right). \end{aligned}$$

Therefore, there exists a mapping  $b_2 : \mathcal{X} \rightarrow \mathbb{R}$  such that

$$\mathbb{E}_{\pi_0} [(R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b_2(x))^2] \leq O\left(\frac{1}{m} e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B \cdot \frac{n+m}{m} \cdot \epsilon_c\right).$$

From Lemma E.4, we have with probability at least  $1 - 2\delta$ , there exists a mapping  $b_1 : \mathcal{X} \rightarrow \mathbb{R}$  such that

$$\mathbb{E}_{\pi_0} [(R(\hat{\theta}_0, x, a) - R(\theta_*, x, a) - b_1(x))^2] \leq O\left(\frac{1}{m} e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B (1 + c_{m,n}) \epsilon_c\right).$$

Hence, we can complete the proof by a union bound over the two events and Assumption 2.6.  $\square$

*Proof of Theorem 4.3.* Let  $b$  be the mapping defined in Lemma E.4 for  $\hat{\theta}$ . We have

$$\begin{aligned} &\mathbb{E}_{\pi_{\theta_*}^\eta} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\theta_*}^\eta(a|x)}{\pi_0(a|x)} \right] - \mathbb{E}_{\pi_{\hat{\theta}}^\eta} \left[ R(\theta_*, x, a) - \frac{1}{\eta} \ln \frac{\pi_{\hat{\theta}}^\eta(a|x)}{\pi_0(a|x)} \right] \\ &= \frac{1}{\eta} \mathbb{E}_{\pi_{\theta_*}^\eta} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\theta_*}^\eta(a|x)} \right] - \frac{1}{\eta} \mathbb{E}_{\pi_{\hat{\theta}}^\eta} \left[ \ln \frac{\pi_0(a|x) \cdot \exp(\eta R(\theta_*, x, a))}{\pi_{\hat{\theta}}^\eta(a|x)} \right] \\ &= \frac{1}{\eta} \mathbb{E}_{x \sim d_0} [\ln Z_{\theta_*}^\eta(x)] - \frac{1}{\eta} \mathbb{E}_{x \sim d_0} [\ln Z_{\hat{\theta}}^\eta(x)] - \mathbb{E}_{x \sim d_0} \left[ \sum_{a \in \mathcal{A}} \pi_{\hat{\theta}}^\eta(a|x) \cdot (R(\theta_*, x, a) - R(\hat{\theta}, x, a)) \right]. \end{aligned}$$

For an arbitrary reward function  $f : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ , let  $\Delta(x, a) = f(x, a) - R(\theta_*, x, a)$ . Consider the following first derivative of  $J(f) = \ln Z_f^\eta(x) - \eta \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \Delta(x, a)$ , where  $Z_f^\eta(x) = \sum_{a \in \mathcal{A}} \pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))$  and  $\pi_f^\eta(a|x) \propto \pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))$ .

Similar to the proof of Theorem 3.3, we still have

$$\begin{aligned} &\frac{\partial}{\partial \Delta(x, a)} \left[ \ln Z_f^\eta(x) - \eta \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \Delta(x, a) \right] \\ &= \frac{1}{Z_f^\eta(x)} \cdot \pi_0(a|x) \exp(\eta \cdot f(x, a)) \cdot \eta - \eta \cdot \pi_f^\eta(a|x) \\ &\quad - \eta \cdot \Delta(x, a) \cdot \frac{\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))}{Z_f^\eta(x)} \cdot \eta + \eta \cdot \Delta(x, a) \cdot \frac{[\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))]^2}{[Z_f^\eta(x)]^2} \cdot \eta \\ &\quad + \eta \sum_{a' \in \mathcal{A} \setminus \{a\}} \frac{\pi_0(a'|x) \cdot \exp(\eta \cdot f(x, a'))}{Z_f^\eta(x)} \cdot \eta \cdot \Delta(x, a') \cdot \frac{\pi_0(a|x) \cdot \exp(\eta \cdot f(x, a))}{Z_f^\eta(x)} \end{aligned}$$

$$= -\eta^2 \pi_f^\eta(a|x) \Delta(x, a) + \eta^2 [\pi_f^\eta(a|x)]^2 \cdot \Delta(x, a) + \eta^2 \sum_{a' \in \mathcal{A} \setminus \{a\}} \pi_f^\eta(a'|x) \pi_f^\eta(a|x) \Delta(x, a').$$

Note that

$$\begin{aligned} J(R(\hat{\theta}, x, \cdot)) &= \ln Z_{\hat{\theta}}^\eta(x) - \eta \sum_{a \in \mathcal{A}} \pi_{\hat{\theta}}^\eta(a|x) \cdot (R(\hat{\theta}, x, a) - R(\theta_*, x, a)) \\ &= \ln \sum_{a \in \mathcal{A}} \pi_0(a|x) \cdot \exp(\eta(R(\hat{\theta}, x, a) - b(x))) - \eta \sum_{a \in \mathcal{A}} \pi_{\hat{\theta}}^\eta(a|x) \cdot (R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b(x)) \\ &= J(R(\hat{\theta}, x, \cdot) - b(x)). \end{aligned}$$

Therefore, there exists  $f(\cdot, \cdot) = \gamma[R(\hat{\theta}, \cdot, \cdot) - b(\cdot)] + (1 - \gamma)R(\theta_*, \cdot, \cdot)$  such that  $(\gamma \in (0, 1))$

$$\begin{aligned} &\mathbb{E}_{x \sim d_0} [J(R(\hat{\theta}, \cdot, \cdot)) - J(R(\theta_*, \cdot, \cdot))] \\ &= \frac{1}{\eta} \mathbb{E}_{x \sim d_0} \left[ -\eta^2 \sum_{a \in \mathcal{A}} \pi_f^\eta(a|x) \cdot \gamma \cdot (R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b(x))^2 \right] \\ &\quad + \frac{1}{\eta} \mathbb{E}_{x \sim d_0} \left[ \gamma \eta^2 \sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \pi_f^\eta(a_1|x) \pi_f^\eta(a_2|x) (R(\hat{\theta}, x, a_1) - R(\theta_*, x, a_1) - b(x)) \right. \\ &\quad \left. (R(\hat{\theta}, x, a_2) - R(\theta_*, x, a_2) - b(x)) \right] \\ &\geq -\eta \cdot \mathbb{E}_{\pi_f^\eta} [(R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b(x))^2] \end{aligned}$$

From Lemma E.2, if  $m \geq 32\eta^2 D^2 e^B \cdot \log(2N_{\mathcal{R}}(\epsilon_c)/\delta)$ , for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$  such that  $\pi_0(a|x) > 0$ , it holds that

$$\eta |R(\hat{\theta}_0, x, a) - R(\theta_*, x, a) - b_1(x)| \leq 1, \quad \eta |R(\hat{\theta}_0, x, a) - R(\theta_*, x, a) - b_2(x)| \leq 1,$$

which means that

$$\frac{\pi_f^\eta}{\pi_{\hat{\theta}_0}^\eta} \leq e^4.$$

Let  $\epsilon_c = \min\left\{\frac{\epsilon}{2(1+c_{m,n})e^B}, \frac{1}{(1+c_{m,n})e^B \eta^2 D^2}\right\}$ . From Lemma E.4, under the condition of the theorem, with high probability the output policy  $\pi_{\hat{\theta}}^\eta$  is  $O(\epsilon)$  optimal.  $\square$

### E.3 PROOF OF THEOREM ??

In this subsection, we also discuss our result under the local coverage condition (Definition 2.8).

**Lemma E.6.** Let  $\hat{\theta}$  be the MLE estimator defined in Algorithm 2. Then for any threshold  $\epsilon_c > 0$ , with probability at least  $1 - 2\delta$ , it holds that

$$\mathbb{E}_{\pi_0} [(R(\hat{\theta}, x, a) - R(\theta_*, x, a) - b(x))^2] \leq O\left(\frac{1}{m} e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B \cdot \frac{m+n}{m} \cdot \epsilon_c\right)$$

for some mapping  $b$  from  $\mathcal{X} \rightarrow \mathbb{R}$ .

*Proof.* From Lemma E.3, with probability at least  $1 - \delta$ ,

$$\begin{aligned} &\sum_{i=1}^m ([R(\hat{\theta}, \tilde{x}_i, \tilde{a}_i^2) - R(\hat{\theta}, \tilde{x}_i, \tilde{a}_i^1)] - [R(\theta_*, \tilde{x}_i, \tilde{a}_i^2) - R(\theta_*, \tilde{x}_i, \tilde{a}_i^1)])^2 \\ &\leq O(e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B(m+n)\epsilon_c). \end{aligned}$$

Following the same argument as (E.4) in the proof of Lemma E.4 and using the union bound, we have with probability at least  $1 - 2\delta$ ,

$$\mathbb{E}_{x \sim d_0} \mathbb{E}_{a^1, a^2 \sim \pi^0} [(R(\hat{\theta}, x, a^1) - R(\theta_*, x, a^1) - R(\hat{\theta}, x, a^2) + R(\theta_*, x, a^2))^2]$$

$$\leq O\left(\frac{1}{m}e^B \log(N_{\mathcal{R}}(\epsilon_c)/\delta) + e^B \cdot \frac{m+n}{m} \cdot \epsilon_c\right).$$

Then we can complete the proof by setting

$$b(x) = \mathbb{E}_{\pi_0}[R(\hat{\theta}, x, a) - R(\theta_*, x, a)].$$

□

## F AUXILIARY LEMMAS

**Lemma F.1** (Freedman inequality). Let  $M, v > 0$  be fixed constants. Let  $\{X_i\}_{i=1}^n$  be a stochastic process,  $\{\mathcal{G}_i\}_i$  be a sequence of  $\sigma$ -fields, and  $X_i$  be  $\mathcal{G}_i$ -measurable, while almost surely

$$\mathbb{E}[X_i|\mathcal{G}_i] = 0, |X_i| \leq M, \text{ and } \sum_{i=1}^n \mathbb{E}[X_i^2|\mathcal{G}_{i-1}] \leq v.$$

Then for any  $\delta > 0$ , with probability at least  $1 - \delta$ , it holds that

$$\sum_{i=1}^n X_i \leq \sqrt{2v \log(1/\delta)} + \frac{2}{3}M \log(1/\delta).$$

**Lemma F.2.** Suppose  $a, b \geq 0$ . If  $x^2 \leq a + b \cdot x$ , then  $x^2 \leq 2b^2 + 2a$ .

*Proof.* By solving the root of quadratic polynomial  $q(x) := x^2 - b \cdot x - a$ , we obtain  $\max\{x_1, x_2\} = (b + \sqrt{b^2 + 4a})/2$ . Hence, we have  $x \leq (b + \sqrt{b^2 + 4a})/2$  provided that  $q(x) \leq 0$ . Then we further have

$$x^2 \leq \frac{1}{4} \left(b + \sqrt{b^2 + 4a}\right)^2 \leq \frac{1}{4} \cdot 2(b^2 + b^2 + 4a) \leq 2b^2 + 2a. \quad (\text{F.1})$$

□

**Lemma F.3** (Pinsker's inequality). If  $\mathbb{P}_1, \mathbb{P}_2$  are two probability measures on a common measurable space  $(\Omega, \mathcal{F})$ , then it holds that

$$\delta(\mathbb{P}_1, \mathbb{P}_2) \leq \sqrt{\frac{1}{2} \text{KL}(\mathbb{P}_1 \| \mathbb{P}_2)},$$

where  $\delta(\cdot, \cdot)$  is the total variation distance and  $\text{KL}(\cdot \| \cdot)$  is the Kullback-Leibler divergence.

**Lemma F.4** (Lemma A.4, Foster et al. 2021). For any sequence of real-valued random variables  $(X_t)_{t \leq T}$  adapted to a filtration  $(\mathcal{F}_t)_{t \leq T}$ , it holds that with probability at least  $1 - \delta$ , for all  $T' \leq T$ ,

$$\sum_{t=1}^{T'} X_t \leq \sum_{t=1}^{T'} \log(\mathbb{E}_{t-1}[e^{X_t}]) + \log(1/\delta).$$