
DRO-REBEL: Distributionally Robust Relative-Reward Regression for Fast and Efficient LLM Alignment

Sharan Sahu

Department of Statistics and Data Science
Cornell University
Ithaca, NY 14853
ss4329@cornell.edu

Abstract

Reinforcement Learning with Human Feedback (RLHF) has become crucial for aligning Large Language Models (LLMs) with human intent. However, existing offline RLHF approaches suffer from overoptimization, where language models degrade by overfitting inaccuracies and drifting from preferred behaviors observed during training [Huang et al., 2025]. Recent methods introduce Distributionally Robust Optimization (DRO) to address robustness under preference shifts, but these methods typically lack sample efficiency, rarely consider diverse human preferences, and use complex heuristics [Mandal et al., 2025, Xu et al., 2025, Chakraborty et al., 2024]. We propose *DRO-REBEL*, a unified family of robust REBEL updates [Gao et al., 2024] instantiated with type- p Wasserstein, Kullback-Leibler (KL), and χ^2 ambiguity sets. Leveraging Fenchel duality, each DRO-REBEL step reduces to a simple relative-reward regression, preserving REBEL’s scalability and avoiding PPO-style clipping or value networks. Under standard linear-reward and log-linear policy classes with a data-coverage assumption, we prove “slow-rate” $O(n^{-1/4})$ estimation-error bounds featuring substantially tighter constants than prior DRO-DPO methods, and we further recover the minimax-optimal “fast-rate” $O(n^{-1/2})$ via a localized Rademacher complexity argument. Crucially, by adapting our localized-complexity analysis to the WDPO and KLDPO algorithms of Xu et al. [2025], we close their existing $O(n^{-1/4})$ vs. $O(n^{-1/2})$ gap and show that both Wasserstein DPO and KL-DPO likewise attain the optimal parametric rate. We also derive tractable SGD-based algorithms for each divergence—using gradient regularization for Wasserstein, importance weighting for KL, and an efficient 1-D dual solve for χ^2 —and demonstrate DRO-REBEL’s superior robustness and sample efficiency on both synthetic benchmarks and language model alignment tasks.

1 Introduction

RLHF has emerged as one of the most important stages of aligning LLMs with human intent [Christiano et al., 2023, Ziegler et al., 2020]. Typically after supervised fine-tuning (SFT), an additional alignment phase is often required to refine their behavior based on human feedback. The alignment of LLMs with human values and preferences is a central objective in machine learning, enabling these models to produce outputs that are useful, safe, and aligned with human intent. In RLHF, human evaluators provide preference rankings that are subsequently utilized to train a reward model, guiding a policy optimization step to maximize learned rewards [Ouyang et al., 2022]. Despite its success, standard RLHF methodologies are fragile mainly due to three reasons: (i) *Assumption that one reward model can model diverse human preferences*: Many RLHF methodologies including popular methods such as Direct Preference Optimization (DPO) [Rafailov et al., 2024] and Proximal Policy Optimization (PPO) [Schulman et al., 2017] assume that a single reward function can model and accurately capture diverse human preferences. In reality, human preferences are highly diverse, context-dependent, and distributional, making it infeasible to represent them within one single reward function. To this end, there has been work done in creating Bayesian frameworks for robust reward modeling [Yan et al., 2024], modeling loss as a weighted combination of different topics and using out-of-distribution detection to reject bad behavior [Bai et al., 2022], or formulating a mixture of reward models [Chakraborty et al., 2024]. (ii) *Reward hacking*: Alignment depends on the quality of the human preference data collected. Unfortunately, this process

is inherently noisy and prone to bias, conflicting opinions, and inconsistency which leads to misaligned preference estimation. This issue is exacerbated by reward hacking where instead of learning reward functions that are aligned with genuine human intent, models learn undesirable shortcuts to maximize the estimated reward function. Subsequently, these models appear to generate responses that appear aligned but deviate from human intent. There are some works that directly address this such as [Bukharin et al., 2024]. (iii) *Distribution shift*: Standard RLHF alignment algorithms use static preference datasets for training, collected under controlled conditions. However, the preferences of real-world users can often be out-of-distribution from that of the training data, depending on several factors such as geographic location, demographics, etc. Thus, a language model in the face of distribution shift may see catastrophic degradation in performance due to overfitting inaccuracies and diverging from human-preferred responses encountered in training data [LeVine et al., 2024, Kirk et al., 2024, Casper et al., 2023]. We focus on the problem of distribution shift, also known as *overoptimization* [Huang et al., 2025].

Recently, distributionally robust RLHF methods have emerged to tackle robustness challenges under distributional shifts in prompts and preferences [Mandal et al., 2025, Xu et al., 2025]. Specifically, Mandal et al. [2025] and Xu et al. [2025] introduced DRO variants of popular RLHF methods, namely DPO and PPO, employing uncertainty sets defined via Chi-Squared (χ^2), type- p Wasserstein, and Kullback–Leibler (KL) divergences. Unfortunately, it is known that PPO requires multiple heuristics to enable stable convergence (e.g. value networks, clipping), and is notorious for its sensitivity to the precise implementation of these components. Recently, Gao et al. [2024] proposed REBEL, an algorithm that cleanly reduces the problem of policy optimization to regressing the relative reward between two completions to a prompt in terms of the policy. They find that REBEL avoids the use of "unjustified" heuristics like PPO and enjoys strong convergence and regret guarantees, similar to Natural Policy Gradient [Kakade, 2001], while also being scalable due to not requiring inversion of the Fisher information matrix. It should be noted that REBEL is much more sample efficient compared to methods like DPO and PPO.

Our contributions. Inspired by the strong theoretical guarantees of REBEL in terms of sample efficiency and simplicity, we introduce *DRO–REBEL*, a family of robust REBEL updates for RLHF under distributional shifts, and make the following advances:

- **Unified robust updates via duality.** We extend REBEL to type- p Wasserstein, KL and χ^2 ambiguity sets, showing that each robust policy update admits a simple relative-reward regression via Fenchel (or strong) duality, preserving REBEL’s scalability and eliminating heuristic tweaks.
- **Sharper “slow-rate” guarantees.** Under standard linear-reward and log-linear-policy assumptions plus a data-coverage condition, we prove an $O(n^{-1/4})$ estimation-error bound for every DRO–REBEL variant. By replacing logistic links with linear regression, we eliminate hidden exponential factors and tighten constants over prior DRO-DPO analyses [Xu et al., 2025, Mandal et al., 2025].
- **Minimax-optimal “fast” rates.** Through a localized Rademacher-complexity argument, we recover the optimal $O(n^{-1/2})$ parametric convergence rate for all DRO–REBEL variants, closing the gap between robust and non-robust RLHF methods.
- **Improved curvature and stability.** We show DRO–REBEL’s strong-convexity modulus scales only with the data-coverage constant λ and step-size η , avoiding the exponential sensitivity to the Bradley–Terry curvature in WDPO/KLDPO [Xu et al., 2025]. This yields uniformly smaller excess-risk bounds and more stable updates.
- **Accelerated robust DPO analysis.** We extend our fast-rate machinery to distributionally robust Direct Preference Optimization (WDPO and KLDPO), proving an $O(n^{-1/2})$ estimation-error rate, improving on the $O(n^{-1/4})$ rates in prior DRO-DPO theory [Xu et al., 2025, Mandal et al., 2025].

1.1 Related Work

Robust RLHF: There has been some recent work in this area that aims to address RLHF overoptimization. Bai et al. [2022] propose addressing distribution shift by adjusting the weights on the combination of loss functions based on different topics (harmless vs. helpful) for robust reward learning. They also propose using out-of-distribution detection to filter and reject known types of bad behavior. [Chakraborty et al., 2024] proposes a MaxMin approach to RLHF, using mixtures of reward models to honor diverse human preference distributions through an expectation-maximization approach, and a robust policy based on these rewards via a max-min optimization. In a similar vein, Padmakumar et al. [2024] tries to augment the human preference datasets with synthetic preference judgments in order to estimate the diversity of user preferences. There has also been some foundational theoretical work towards this problem. Yan et al. [2024] proposed a Bayesian reward model ensemble to model the uncertainty set of the reward functions and systematically choose rewards in the uncertainty set with the tightest confidence band. Another line of work focuses on robust reward modeling as an alternative to distributionally robust optimization. For instance, Bukharin et al. [2024] propose R3M, a method that explicitly models corrupted preference labels as sparse outliers. They formulate reward

learning as an ℓ_1 -regularized maximum likelihood estimation problem, enabling robust recovery of the underlying reward function even in the presence of noisy or inconsistent human feedback. While our work focuses on embedding robustness at the policy optimization level using distributional uncertainty sets (e.g., χ^2 , and Wasserstein), R3M represents a complementary direction that enhances robustness by improving the reliability of the reward model itself.

Robust DPO: There have been several works that approach this problem using DRO. Huang et al. [2025] proposed χ PO that implements the principle of pessimism in the face of uncertainty via regularization with the χ^2 -divergence for avoiding reward hacking/overoptimization with respect to the estimated reward. Wu et al. [2024] focus on noisy preference data and categorize the types of noise in DPO, introducing Dr. DPO to improve pairwise robustness through a DRO formulation with a tunable reliability parameter. Hong et al. [2024] propose an adaptive preference loss grounded in DRO that adjusts scaling weights across preference pairs to account for ambiguity in human feedback, enhancing reward estimation flexibility and policy performance. Separately, Zhang et al. [2024] introduce a lightweight uncertainty-aware approach called AdvPO, combining last-layer embedding-based uncertainty estimation with a DRO formulation to address overoptimization in reward-based RLHF. There are two related works that are most similar to our approach. Xu et al. [2025] develop Wasserstein and KL-based DRO formulations of Direct Preference Optimization (WDPO and KLDPO), providing sample complexity bounds and scalable gradient-based algorithms. Their methods achieve improved alignment performance under shifting user preference distributions. Similarly, Mandal et al. [2025] propose robust variants of both reward-based and reward-free RLHF methods, incorporating DRO into the reward estimation and policy optimization phases using Total Variation and Wasserstein distances. Their algorithms retain the structure of existing RLHF pipelines while providing theoretical convergence guarantees and demonstrating robustness to out-of-distribution (OOD) tasks.

Distributionally Robust Learning: The DRO framework has been applied to various areas ranging from supervised learning [Namkoong and Duchi, 2017b, Shah et al., 2020], reinforcement learning [Zhang et al., 2020, Yang et al., 2021], and multi-armed bandits [Gao et al., 2022, Zhou et al., 2022]. There is a wealth of theoretical results using f-divergences and Wasserstein distances developed for tackling problems in this setting [Duchi and Namkoong, 2022, Shapiro and Xu, 2022].

2 Preliminaries

2.1 Notations

We will denote sets using calligraphic letters i.e. $\mathcal{S}, \mathcal{A}, \mathcal{Z}$. When we refer to the norm $\|x\|$, we are referring to the Euclidean norm. For a measure \mathbb{P} , we refer to the empirical measure \mathbb{P}_n to mean drawing samples $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}$ with $\mathbb{P}_n = 1/n \sum_{i=1}^n \delta_{x_i}$ where δ is the Dirac measure. We denote $\ell(z; \theta)$ to be the loss incurred by sample z with policy parameter θ . We denote $\mathcal{M}(\mathcal{Z})$ to be the set of Borel measures supported on set \mathcal{Z} . Lastly, we denote $\lambda_{\min}(A)$ to be the minimum eigenvalue of a symmetric matrix $A \in \mathbb{S}^n$.

2.2 Divergences

In this section, we will define the divergences that we will use to define our ambiguity sets in the DRO setting.

Definition 2.1 (Type-p Wasserstein Distance). *The type-p ($p \in [1, \infty)$) Wasserstein distance between two distributions $\mathbb{P}, \mathbb{Q} \in \mathcal{M}(\Xi)$ is defined as*

$$\mathcal{W}_p(\mathbb{P}, \mathbb{Q}) = \left(\inf_{\pi \in \Pi(\mathbb{P}, \mathbb{Q})} \int_{\mathbb{R}^d \times \mathbb{R}^d} d(\xi, \eta)^p \pi(d\xi, d\eta) \right)^{1/p}$$

where π is a coupling between the marginal distributions $\xi \sim \mathbb{P}$ and $\eta \sim \mathbb{Q}$ and d is a pseudometric defined on \mathcal{Z} .

Definition 2.2 (Kullback-Leibler (KL) Divergence).

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) = \int_{\Xi} \log \left(\frac{d\mathbb{P}}{d\mathbb{Q}} \right) d\mathbb{P}$$

Definition 2.3 (Chi-Squared Divergence). *If $\mathbb{P} \ll \mathbb{Q}$,*

$$D_{\chi^2}(\mathbb{P} \parallel \mathbb{Q}) = \int_{\Xi} \left(\frac{d\mathbb{P}}{d\mathbb{Q}} - 1 \right)^2 d\mathbb{Q}$$

Using these, we can define our ambiguity sets as follows

Definition 2.4 (Distributional Uncertainty Sets). Let $\varepsilon > 0$ and $\mathbb{P}^\circ \in \mathcal{M}(\mathcal{Z})$. Then, we define the ambiguity set as

$$\mathcal{B}_\varepsilon(\mathbb{P}^\circ; D) = \{\mathbb{P} \in \mathcal{M}(\mathcal{Z}) : D(\mathbb{P}, \mathbb{P}^\circ) \leq \varepsilon\}$$

where $D(\cdot, \cdot)$ is a distance metric between two probability distributions i.e. type- p Wasserstein, KL, χ^2 .

2.3 Reinforcement Learning from Human Feedback (RLHF)

Reinforcement Learning from Human Feedback (RLHF), as introduced by Christiano et al. [2023] and later adapted by Ouyang et al. [2022], consists of two primary stages: (1) learning a reward model from preference data, and (2) optimizing a policy to maximize a KL-regularized value function. We assume access to a preference dataset $\mathcal{D}_{\text{src}} = \{(x, a^1, a^2)\}$, where $x \in \mathcal{S}$ is a prompt, and $a^1, a^2 \in \mathcal{A}$ are two possible completions of the prompt x generated from a reference policy $\pi_{\text{SFT}}(\cdot | x)$ (e.g., a supervised fine-tuned model). $\pi_{\text{SFT}}(\cdot | x)$ involves fine-tuning a pre-trained LLM through supervised learning on high-quality data, curated for downstream tasks. A human annotator provides preference feedback indicating $a^1 \succ a^2 | x$. The most common model for preference learning is the Bradley-Terry (BT) model, which assumes that

$$\begin{aligned} \mathcal{P}^*(a^1 \succ a^2 | x) &= \sigma(r^*(x, a^1) - r^*(x, a^2)) \\ &= \frac{1}{1 + \exp(r^*(x, a^1) - r^*(x, a^2))}, \end{aligned}$$

where r^* is the underlying (unknown) reward function used by the annotator. The first step in RLHF is to learn a parametrized reward model $r_\phi(s, a)$ by solving the following maximum likelihood estimation problem:

$$r_\phi \leftarrow \arg \min_{r_\phi} -\mathbb{E}_{(x, a^1, a^2) \sim \mathcal{D}_{\text{src}}} [\log \sigma(r_\phi(x, a^1) - r_\phi(x, a^2))].$$

Given the learned reward model r_ϕ , the second step is to solve the KL-regularized policy optimization problem:

$$\pi_\theta \leftarrow \arg \max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}_{\text{src}}, y \sim \pi_\theta(\cdot | x)} \left[r_\phi(x, y) - \beta \log \frac{\pi_\theta(y | x)}{\pi_{\text{SFT}}(y | x)} \right],$$

where β controls the deviation between the learned and reference policy.

2.4 REBEL: Regression-Based Policy Optimization

Let (x, a) represent a *prompt-response* pair, where $x \in \mathcal{S}$ is a context or prompt, and $a \in \mathcal{A}$ is a response (e.g., a sequence of tokens or actions). We assume access to a reward function $r(x, a)$, which may be a learned preference model [Christiano et al., 2023]. Let $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ be a stochastic policy mapping prompts to distributions over responses. We denote the prompt distribution as ρ , and let $\pi_\theta(a | x)$ denote a parameterized policy with parameters θ . The **REBEL** (REGression to RELative REward-Based RL) [Gao et al., 2024] algorithm directly regresses to relative reward differences through KL-constrained updates. The REBEL algorithm is detailed in Algorithm 1.

Algorithm 1 REBEL: REGression to RELative REward-Based RL

- 1: **Input:** Reward function r , policy class $\Pi = \{\pi_\theta : \theta \in \Theta\}$, base distribution μ , learning rate η
- 2: Initialize policy π_{θ_0}
- 3: **for** $t = 0$ to $T - 1$ **do**
- 4: Collect dataset $\mathcal{D}_t = \{(x, a^1, a^2)\}$ with $x \sim \rho$, $a^1 \sim \pi_t(\cdot | x)$, $a^2 \sim \mu(\cdot | x)$
- 5: Update policy by solving:

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \sum_{(x, a^1, a^2) \in \mathcal{D}_t} \left(\frac{1}{\eta} \left[\log \frac{\pi_\theta(a^1 | x)}{\pi_{\theta_t}(a^1 | x)} - \log \frac{\pi_\theta(a^2 | x)}{\pi_{\theta_t}(a^2 | x)} \right] - [r(x, a^1) - r(x, a^2)] \right)^2 \quad (1)$$

6: **end for**

At each iteration, REBEL approximates the solution to a KL-constrained policy optimization objective:

$$\pi_{t+1} = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \rho, a \sim \pi(\cdot | x)} [r(x, a)] - \frac{1}{\eta} \mathbb{E}_{x \sim \rho} [\text{KL}(\pi(\cdot | x) \parallel \pi_t(\cdot | x))],$$

which encourages reward maximization while regularizing the policy to remain close to the previous iterate. This objective is particularly well-suited for fine-tuning language models using learned or noisy reward signals while maintaining stability. Adapting REBEL for distributionally robust RLHF is particularly appealing because it offers both theoretical and practical advantages over existing methods like PPO and DPO. Whereas PPO relies on heuristic mechanisms (e.g., clipping, value baselines) and DPO requires strong assumptions about preference modeling, REBEL reduces policy optimization to a sequence of regression problems on relative rewards—eliminating the need for explicit value functions or constrained optimization. This simplicity translates into significantly lower sample complexity. In particular, REBEL can achieve convergence guarantees comparable to or better than NPG, with a sample complexity that scales favorably due to its variance-reduced gradient structure. Empirically, REBEL has been shown to converge faster than PPO and outperform DPO in both language and image generation tasks. Building on this regression-based perspective, our DRO-REBEL algorithms simply replace the standard squared-error loss in each REBEL update with its robust counterpart under the chosen divergence (Wasserstein, KL, or χ^2). As a result, DRO-REBEL inherits REBEL’s stability and low sample complexity while gaining worst-case robustness guarantees under distributional shifts.

2.5 Direct Preference Optimization (DPO)

The DPO [Rafailov et al., 2024] procedure is a form of offline RLHF which avoids issues in previous policy optimization algorithms like PPO [Schulman et al., 2017] by identifying a mapping (and re-parametrization) between language model policies and reward functions that enables training LMs to satisfy human preferences directly without using RL or even doing reward model fitting. That is, DPO makes use of the following policy objective

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{SFT}}) = -\mathbb{E}_{(x, a^1, a^2) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(a^1 | x)}{\pi_{\text{SFT}}(a^1 | x)} - \beta \log \frac{\pi_\theta(a^2 | x)}{\pi_{\text{SFT}}(a^2 | x)} \right) \right] \quad (2)$$

One can arrive at this policy objective by observing that the objective in the RL fine-tuning phase has a closed form solution of the form

$$\pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{SFT}}(y | x) \exp \left(\frac{1}{\beta} r(x, y) \right)$$

where $Z(x) = \sum_y \pi_{\text{SFT}}(y | x) \exp \left(\frac{1}{\beta} r(x, y) \right)$. Taking logs and moving terms around, we get

$$r(x, y) = \beta \frac{\pi(y | x)}{\pi_{\text{SFT}}(y | x)} + \beta \log Z(x)$$

Using the BT model and plugging in this re-parametrization, we get the DPO policy objective. One might note that this is not a proper re-parametrization since $Z(x)$ is dependent on r . However, it turns out that defining this re-parametrization as a function from an equivalence class of reward functions to a particular policy is well-defined, does not constrain the class of learned reward models, and allows for the exact recovery of the optimal policy.

3 Distributionally Robust REBEL and DPO

In this section, we will formally define the DRO variants of REBEL and DPO under type-p Wasserstein, KL, and χ^2 divergence ambiguity sets. We must first define the nominal data-generating distribution. Our definitions follow those stated by Xu et al. [2025]. Recall the sampling procedure mentioned in Section 2.3: We have some initial prompt $x \in \mathcal{S}$ that we will assume is sampled from some prompt distribution ρ . We will sample two responses $a^1, a^2 \stackrel{\text{i.i.d.}}{\sim} \pi_{\text{SFT}}(\cdot | x)$. Following Zhu et al. [2023], let $y \in \{0, 1\}$ be a Bernoulli random variable where $y = 1$ if $a^1 \succ a^2 | x$ and $y = 0$ if $a^2 \succ a^1 | x$ with probability corresponding to the Bradley-Terry model \mathcal{P}^* . Using this, we can now define the nominal data-generating distribution.

Definition 3.1 (Nominal Data-Generating Distribution). *Let $\mathcal{Z} = \mathcal{S} \times \mathcal{A} \times \mathcal{A} \times \{0, 1\}$. Then, we define the nominal data-generating distribution as follows*

$$\mathbb{P}^\circ(x, a^1, a^2, y) = \rho(x) \pi_{\text{SFT}}(a^1 | x) \pi_{\text{SFT}}(a^2 | x) [\mathbb{1}_{\{y=1\}} \mathcal{P}^*(a^1 \succ a^2 | x) + \mathbb{1}_{\{y=0\}} \mathcal{P}^*(a^2 \succ a^1 | x)]$$

where $x \sim \rho$ and $y \sim \text{Ber}(\mathcal{P}^*(a^1 \succ a^2 | \cdot))$. We will denote $z = (x, a^1, a^2, y) \in \mathcal{Z}$ and $\mathbb{P}^\circ(z) = \mathbb{P}^\circ(x, a^1, a^2, y)$. We will also assume \mathbb{P}° generates dataset $\mathcal{D} = \{z_i\}_{i=1}^n$ used for learning i.e. $z_i \sim \mathbb{P}^\circ$.

3.1 Distributionally Robust REBEL

From the REBEL update (Equation (1)), we define the pointwise loss as follows

$$\ell_{\text{REBEL}}(z; \theta) = \left(\frac{1}{\eta} \left[\log \frac{\pi_{\theta}(a^1 | x)}{\pi_t(a^1 | x)} - \log \frac{\pi_{\theta}(a^2 | x)}{\pi_t(a^2 | x)} \right] - [r(x, a^1) - r(x, a^2)] \right)^2$$

For $\varepsilon > 0$, define the ambiguity set as $\mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)$ for nominal distribution \mathbb{P}° and distance measure D . Using the DRO framework, we consider the following distributionally robust optimization problem:

$$\min_{\theta} \max_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)} \mathbb{E}_{z \sim \mathbb{P}} [\ell_{\text{REBEL}}(z; \theta)]$$

which directly captures our objective: finding the best policy under worst-case distributional shift. Now, let us define the following D -DRO-REBEL loss function:

$$\mathcal{L}^D(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)} \mathbb{E}_{z \sim \mathbb{P}} [\ell_{\text{REBEL}}(z; \theta)]$$

where $\mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)$ denotes an ambiguity set centered at the nominal distribution \mathbb{P}° , defined using a divergence or distance D . This formulation is general and allows us to instantiate a family of distributionally robust REBEL objectives by choosing different D —such as the type- p Wasserstein distance, Kullback–Leibler (KL) divergence, or chi-squared (χ^2) divergence. Each choice of D yields a different robustness profile and tractable dual formulation, enabling us to tailor the algorithm to specific distributional shift scenarios. When the nominal distribution \mathbb{P}° is replaced with its empirical counterpart, i.e., $\mathbb{P}_n^{\circ} := \frac{1}{n} \sum_{i=1}^n \delta_{z_i}$, where z_1, \dots, z_n are n i.i.d. samples from \mathbb{P}° , we use $\mathcal{L}_n^D(\theta; \varepsilon)$ to denote the empirical D -REBEL loss incurred by the policy parameter θ .

3.2 Distributionally Robust DPO

Similar to DRO-REBEL, we can define a distributionally robust counterpart for DPO. From the DPO loss function (Equation 2.5), we define the following pointwise loss:

$$\ell_{\text{DPO}}(z; \theta) = -y \log \sigma(\beta h_{\theta}) - (1 - y) \log \sigma(-\beta h_{\theta})$$

where $h_{\theta}(s, a^1, a^2) := \log \frac{\pi_{\theta}(a^1 | s)}{\pi_{\text{ref}}(a^1 | s)} - \log \frac{\pi_{\theta}(a^2 | s)}{\pi_{\text{ref}}(a^2 | s)}$ is the preference score of an answer a^1 relative to answer a^2 . As we did for REBEL, we can use the DRO framework to formulate the following distributionally robust optimization problem:

$$\min_{\theta} \max_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)} \mathbb{E}_{z \sim \mathbb{P}} [\ell_{\text{DPO}}(z; \theta)]$$

We define the following D-DPO loss function:

$$\mathcal{L}^D(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^{\circ}; D)} \mathbb{E}_{z \sim \mathbb{P}} [\ell_{\text{DPO}}(z; \theta)]$$

4 Theoretical Results

In this section, we will provide several sample complexity results for DRO-REBEL under the ambiguity sets previously mentioned. First, we state some assumptions that we make in our analysis

Assumption 1 (Linear reward class). *Let $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a known d -dimensional feature mapping with $\sup_{x,a} \|\phi(x, a)\|_2 \leq 1$ and $\omega \in \mathbb{R}^d$ such that $\|\omega\|_2 \leq F$ for $F > 0$. We consider the following class of linear reward functions:*

$$\mathcal{F} := \{r_{\omega} : r_{\omega}(x, a) = \phi(x, a)^{\top} \omega\}$$

Remark 1. *These are standard assumptions in the theoretical analysis of RL algorithms [Agarwal et al., 2021b, Modi et al., 2020], RLHF [Zhu et al., 2023], and DPO [Nika et al., 2024, Chowdhury et al., 2024]. Our analysis can be extended to settings involving neural network-parameterized policy or reward classes, commonly found in modern RLHF for LLM alignment. In such cases, the key properties of the loss function $\ell(z; \theta)$ (e.g., uniform boundedness,*

Lipschitz continuity, and strong convexity with respect to θ) would require corresponding assumptions on the neural network architectures, such as twice differentiability, smoothness (bounded Hessian), and specific architectural choices (e.g., linear output layers for strong convexity) that ensure these properties hold for the composite function.

Assumption 2 (Log-linear policy class). Let $\psi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a known d -dimensional feature mapping with $\max_{x,a} \|\psi(x,a)\|_2 \leq 1$. Assume a bounded policy parameter set $\Theta := \{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq B\}$. We consider the following class of log-linear policies:

$$\Pi := \left\{ \pi_\theta : \pi_\theta(a | x) = \frac{\exp(\theta^\top \psi(x,a))}{\sum_{a' \in \mathcal{A}} \exp(\theta^\top \psi(x,a'))} \right\}.$$

Remark 2. These are standard assumptions in the theoretical analysis of RL algorithms [Agarwal et al., 2021b, Modi et al., 2020], RLHF [Zhu et al., 2023], and DPO [Nika et al., 2024, Chowdhury et al., 2024]. Our analysis can be extended to neural policy classes where $\theta^\top \psi(s,a)$ is replaced by $f_\theta(s,a)$, where f_θ is a neural network satisfying twice differentiability and smoothness assumptions.

We also make the following data coverage assumption on the uncertainty set $\mathcal{B}_\varepsilon(\mathbb{P}^\circ; D)$:

Assumption 3 (Regularity condition). There exists $\lambda > 0$ such that

$$\Sigma_{\mathbb{P}} := \mathbb{E}_{(x,a^1,a^2,y) \sim \mathbb{P}} \left[(\psi(x,a^1) - \psi(x,a^2)) (\psi(x,a^1) - \psi(x,a^2))^\top \right] \succeq \lambda I, \forall \mathbb{P} \in \mathcal{P}(\rho; \mathbb{P}^\circ).$$

Remark 3. Similar assumptions on data coverage under linear architecture models are standard in the offline RL literature [Agarwal et al., 2021a, Wang et al., 2020, Jin et al., 2022]. Implicitly, Assumption 2 imposes $\lambda \leq \lambda_{\min}(\Sigma_{\mathbb{P}^\circ})$, meaning the data-generating distribution \mathbb{P}° must have sufficient coverage.

4.1 "Slow Rate" Estimation Errors

Define $\theta^{\mathcal{W}_p} \in \arg \min_{\theta \in \Theta} \mathcal{L}^{\mathcal{W}_p}(\theta)$ be the true optimal policy estimate and the empirical estimate as $\hat{\theta}_n^{\mathcal{W}_p} \in \arg \min_{\theta \in \Theta} \mathcal{L}_n^{\mathcal{W}_p}(\theta)$. First, we provide a sample complexity result for convergence of robust policy estimation using REBEL. Our proof technique hinges on showing that $\mathcal{L}^{\mathcal{W}_p}$ is strongly convex.

Lemma 1 (Strong convexity of $\mathcal{L}^{\mathcal{W}_p}$). Let $\ell(z; \theta)$ be as defined in the REBEL update. The Wasserstein-DRO-REBEL loss

$$\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell_{\text{REBEL}}(z; \theta)],$$

is $2\lambda/\eta^2$ -strongly convex where λ is from the regularity condition in Assumption 3 and η is from the step size defined in the DRO update 1

We now present our "slow rate" results on the sample complexity for the convergence of the robust policy parameter.

Theorem 1 ("Slow" Estimation error of $\theta^{\mathcal{W}_p}$). Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$

$$\|\theta^{\mathcal{W}_p} - \hat{\theta}_n^{\mathcal{W}_p}\|_2^2 \leq \frac{\eta^2 K_g^2}{\lambda} \sqrt{\frac{2 \log(2/\delta)}{n}}$$

where λ is from the regularity condition in Assumption 3 and $K_g = 4B/\eta + 2F$ where B is from the assumption that the policy parameter set is bounded in Assumption 2 and F is from the Assumption 1.

Proof sketch. For the full proof, we refer readers to Appendix C. At a high level, we first prove that $\ell(z; \theta)$ is uniformly bounded and is $4K_g/\eta$ -Lipschitz in θ where $K_g = 4B/\eta + 2F$. Using this, we can prove that $\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$ is $2/\eta$ -strongly convex in $\|\cdot\|_{\Sigma_{\mathbb{P}}}$. Intuitively taking the supremum over \mathbb{P} should preserve the convex combination and the negative quadratic term and doing this analysis formally allows us to show that $\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon)$ is $2\lambda/\eta^2$ -strongly convex in $\|\cdot\|_2$. The detailed proof for strong convexity can be found in Lemma 19.

Strong duality of Wasserstein DRO [Gao and Kleywegt, 2022] (Corollary 2) allows us to reduce the difference $|\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon)|$ to the concentration $|\mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)]|$ where $\ell_\Delta(z; \theta) = \inf_{z' \in \mathcal{Z}} \{\Delta d^p(z, z') - \ell(z'; \theta)\}$ is the Moreau envelope of $-\ell$. We then use Hoeffding's inequality to obtain a concentration result which is uniform over $\theta \in \Theta$ and Δ . We can now do a "three-term" decomposition of $\mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon)$ into $\mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon)$, $\mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon)$, and $\mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon)$ and bound the first and last term by Hoeffding and the second term by 0. Using strong convexity of $\mathcal{L}^{\mathcal{W}_p}$, we can get the estimation error. The detailed proof for the "slow rate" estimation error can be found at C \square

We prove similar results for KL and χ^2 ambiguity sets using the same ideas used in the Wasserstein ambiguity set setting. We state the "slow rate" estimation rates below and defer the proofs to Appendix D and Appendix E

Theorem 2 ("Slow" Estimation error of θ^{KL}). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$*

$$\|\theta^{\text{KL}} - \hat{\theta}_n^{\text{KL}}\|_2^2 \leq \frac{\eta^2}{\lambda} \sqrt{\frac{2\bar{\lambda}^2 \exp(K_g^2/\bar{\lambda}) \log(2/\delta)}{n}}$$

where $K_g = 4B/\eta + 2F$, $\bar{\lambda}$ is from Assumption 4, and η is the stepsize defined in in the DRO update 1

Theorem 3 ("Slow" Estimation error of θ^{χ^2}). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$*

$$\|\theta^{\chi^2} - \hat{\theta}_n^{\chi^2}\|_2^2 \leq \frac{\eta^2 K_g^2}{\lambda} (1 + K_g^2/4\bar{\lambda}) \sqrt{\frac{2 \log(2/\delta)}{n}}$$

where λ is from the regularity condition in Assumption 3 and $K_g = 4B/\eta + 2F$ where B is from the assumption that the policy parameter set is bounded in Assumption 2, F is from the Assumption 1, and η is from the step size defined in the DRO update 1.

Remark 4 (Estimation rate and improved coverage dependence). *Although both WDPO Xu et al. [2025] and DRO-REBEL achieve the same $n^{-1/4}$ estimation-error rate, DRO-REBEL features substantially tighter constants thanks to its regression-to-relative-rewards formulation and cleaner strong-convexity analysis. In WDPO, the strong-convexity modulus is the product of the Bradley–Terry curvature*

$$\gamma = \frac{\beta^2 e^{4\beta B}}{(1 + e^{4\beta B})^2} \quad \text{and} \quad \lambda,$$

so that the squared-error bound scales as $O(1/(\gamma\lambda))$ and is exponentially sensitive to the logistic scale β [Xu et al., 2025, Lemma 1, Theorem 1]. By contrast, Lemma 19 shows that the DRO-REBEL loss $\mathcal{L}^{\mathcal{W}_p}$ is $(2\lambda/\eta^2)$ -strongly convex, yielding a bound of order $O(1/\lambda)$ (up to the step-size η) with no hidden logistic factors. Concretely, this sharper constant means that for any fixed coverage λ , our excess-risk bound is smaller by the factor $\gamma^{-1} = (1 + e^{4\beta B})^2/(\beta^2 e^{4\beta B})$, which can be enormous when βB is large or preferences are near-degenerate. Moreover, the same phenomenon appears in the KL-DRO setting. Theorem 2 shows

$$\|\theta^{\text{KL}} - \hat{\theta}_n^{\text{KL}}\|_2^2 \leq \frac{\eta^2}{\lambda} \sqrt{\frac{2\bar{\lambda}^2 \exp(K_g^2/\bar{\lambda}) \log(2/\delta)}{n}},$$

where $\bar{\lambda}$ bounds the dual multiplier and L bounds the loss. In prior KLDPO analyses, the dependence on $\exp(K_g^2/\bar{\lambda})$ is tangled with additional Bradley–Terry curvature terms; in DRO-REBEL it appears only through the divergence parameter. Crucially, both Wasserstein and KL results rest on the very same modelling assumptions: linear reward class (Assumption 1), log-linear policy class (Assumption 2), and data-coverage regularity (Assumption 3). This shows that DRO-REBEL’s improvements arise purely from algorithmic simplicity rather than stronger distributional or curvature requirements.

With the above analysis, building on and refining the techniques of Xu et al. [2025], we recover the same $O(n^{-1/4})$ estimation-error rate but with substantially sharper constants. Xu et al. [2025] observe that WDPO’s estimation error decays at $O(n^{-1/4})$, while non-robust DPO already achieves $O(n^{-1/2})$. This slowdown arises because, in the robust setting, one cannot exchange the supremum over distributions with the gradient operator on the empirical robust loss. As a result, the closed-form concentration argument—key to the non-robust analysis—fails. Closing this gap and restoring the optimal inverse square root rate for robust DPO remains an open problem. To close this gap, we develop a *localized Rademacher complexity* analysis for DRO-REBEL which recovers the optimal $O(n^{-1/2})$ convergence rate even under Wasserstein ambiguity.

4.2 A "Master Theorem" for "Fast" Estimation Rates

To state our main guarantee in its cleanest form, we observe that nothing exotic is needed beyond (i) the population DRO objective has a uniform quadratic growth (via our data-coverage assumption), (ii) each per-sample loss $\ell(z; \theta)$ is Lipschitz in θ , which drives the Rademacher/Dudley control, and (iii) each divergence admits a simple Fenchel-duality or direct bound showing $|\mathcal{L}^D(\theta) - \mathbb{E}_P[\ell(z; \theta)]| \leq O(\Delta_n)$ where $\Delta_n = O(n^{-1})$. The following single “master theorem” then automatically yields the parametric $n^{-1/2}$ -rate for all of our DRO variants, Wasserstein, KL, and χ^2 , by plugging in the corresponding Δ_n .

Theorem 4 (Parametric $n^{-1/2}$ -rate for DRO-REBEL). *Let the empirical and population minimizers be defined as:*

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} \mathcal{L}_n^D(\theta; \varepsilon_n), \quad \theta^* = \arg \min_{\theta \in \Theta} \mathcal{L}^D(\theta; \varepsilon_n)$$

Assume the following conditions hold:

1. **Local Strong Convexity:** The population DRO loss $\mathcal{L}^D(\theta; \varepsilon_n)$ is α -strongly convex in a neighborhood of θ^* .
2. **Lipschitz Loss:** The pointwise loss function $\ell(z; \theta)$ is L_g -Lipschitz with respect to θ for all z .
3. **Bounded Loss:** The pointwise loss $\ell(z; \theta)$ is uniformly bounded by K_ℓ for all $z \in \mathcal{Z}, \theta \in \Theta$.
4. **Dual Remainder Bound:** There exists a non-negative quantity Δ_n such that for any $\theta \in \Theta$:

$$|\mathcal{L}^D(\theta; \varepsilon_n) - \mathbb{E}_{\mathbb{P}^\circ}[\ell(z; \theta)]| \leq \Delta_n$$

and

$$|\mathcal{L}_n^D(\theta; \varepsilon_n) - \mathbb{E}_{\mathbb{P}_n^\circ}[\ell(z; \theta)]| \leq \Delta_n$$

Then for any fixed $\delta \in (0, 1)$, with probability at least $1 - \delta$, the estimation error is bounded as:

$$\|\hat{\theta}_n - \theta^*\|_2 \leq \frac{16c_0 L_g}{\alpha} \sqrt{\frac{d}{n}} + \frac{8c_1 L_g}{\alpha} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8\Delta_n}{\alpha}} + \sqrt{\frac{4c_2 K_\ell \log(1/\delta)}{\alpha n}}$$

where the universal constants c_0, c_1, c_2 arise from concentration inequalities. In particular, if the dual remainder bound has rate $\Delta_n = O(n^{-1})$, then the estimator achieves a parametric rate of:

$$\|\hat{\theta}_n - \theta^*\|_2 = O(n^{-1/2})$$

Proof Sketch of Theorem 4. For the full proof, we refer readers to Appendix F. At a high level, for fixed $\theta \in \Theta$, by using the triangle-inequality and the dual-remainder bound, we get $|\mathcal{L}^D(\theta) - \mathcal{L}_n^D(\theta)| \leq 2\Delta_n + |M(\theta) - M_n(\theta)|$. Next, let $r = \|\theta - \theta^*\|_2$ and define $\mathcal{F}_r = \{\ell(\cdot, \theta) - \ell(\cdot, \theta^*) : \|\theta - \theta^*\| \leq r\}$. By a standard symmetrization (Lemma 8), Dudley entropy-integral (Corollary 5), and Bousquet's inequality (Theorem 7) argument, one shows with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}_r} |P_n f - P f| = O(\Delta_n + L_g r \sqrt{(d + \ln(1/\delta))/n} + K_\ell \ln(1/\delta)/n)$$

Let $\mathcal{E}_k = \left\{ \sup_{\|\theta - \theta^*\|_2 \leq r_k} |\mathcal{L}_n^D(\theta; \varepsilon_n) - \mathcal{L}^D(\theta; \varepsilon_n)| \leq \psi_n(r_k, \delta_k) \right\}$ where

$$\psi_n(r, \delta) = O(\Delta_n + L_g r \sqrt{(d + \ln(1/\delta))/n} + K_\ell \ln(1/\delta)/n)$$

We can relate this uniform bound to the estimation error $\tilde{r} = \|\hat{\theta}_n - \theta^*\|_2$ using a **localization argument**. We define a dyadic sequence of radii $r_k = 2^k r_0$ for $k = 0, 1, \dots, \lceil \log_2(R_{\max}/r_0) \rceil$. Using a union bound over these radii, we establish that the uniform convergence event holds simultaneously for all balls $B(\theta^*, r_k)$ with high probability (e.g., $1 - \delta$). Let k^* be the smallest integer such that $\tilde{r} \leq r_{k^*}$. By this very definition of k^* , $\hat{\theta}_n \in B(\theta^*, r_{k^*})$. Furthermore, if $k^* > 0$, then $r_{k^*-1} < \tilde{r} \leq r_{k^*}$. Since $r_{k^*} = 2 \cdot r_{k^*-1}$ (for a dyadic sequence), this implies $r_{k^*} < 2\tilde{r}$. For $k^* = 0$, $\tilde{r} \leq r_0$. In general, we can use $r_{k^*} \leq 2 \max(\tilde{r}, r_0)$ for the linear term of ψ_n , which simplifies to $2\tilde{r}$ when \tilde{r} is large enough (as expected asymptotically). Combining this with the uniform convergence bound for r_{k^*} , α -strong convexity of $\mathcal{L}^D(\theta)$ around θ^* , and the definition of the empirical minimizer, we arrive at a quadratic inequality in \tilde{r} :

$$\frac{\alpha}{2} \tilde{r}^2 - \left(4c_0 L_g \sqrt{\frac{d}{n}} + 2c_1 L_g \sqrt{\frac{\log(1/\delta)}{n}} \right) \tilde{r} - \left(4\Delta_n + 2c_2 \frac{K_\ell}{n} \log(1/\delta) \right) \leq 0$$

Solving this quadratic inequality for \tilde{r} yields

$$\tilde{r} = O\left(\frac{L_g}{\alpha} \sqrt{\frac{d}{n}} + \frac{L_g}{\alpha} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{\Delta_n}{\alpha}} + \sqrt{\frac{K_\ell \log(1/\delta)}{\alpha n}} \right)$$

Taking $\Delta_n = O(n^{-1})$, all terms in the bound become $O(n^{-1/2})$, which establishes the claimed convergence rate. \square

4.3 DRO-REBEL "Fast Rates"

Using this theorem, we can now state the following "fast rate" REBEL estimation results for Wasserstein, KL, and χ^2 . We defer the proofs to Appendix G, H, and I respectively.

Corollary 1 ("Fast" Estimation error of $\theta^{\mathcal{W}_p}$). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$*

$$\|\hat{\theta}_n^{\mathcal{W}_p} - \theta^{\mathcal{W}_p}\|_2 \leq \frac{16c_0K_g\eta}{\lambda} \sqrt{\frac{d}{n}} + \left(\frac{8c_1K_g\eta}{\lambda} + \frac{2K_g\eta\sqrt{c_2}}{\sqrt{\lambda}} \right) \sqrt{\frac{\log(1/\delta)}{n}} + \frac{2\eta\sqrt{L_{\ell,z}}}{\sqrt{\lambda n}}$$

where $c_0, c_1, c_2 > 0$ are some absolute constants, $L_{\ell,z}$ is from Assumption 5, λ is from the regularity condition in Assumption 3, and $K_g = 4B/\eta + 2F$ where B is from the assumption that the policy parameter set is bounded in Assumption 2, and F is from the Assumption 1.

Corollary 2 ("Fast" Estimation error of θ^{KL}). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$*

$$\|\hat{\theta}_n^{\text{KL}} - \theta^{\text{KL}}\|_2 \leq \frac{16c_0K_g\eta}{\lambda} \sqrt{\frac{d}{n}} + \left(\frac{8c_1K_g\eta}{\lambda} + \frac{2K_g\eta\sqrt{c_2}}{\sqrt{\lambda}} \right) \sqrt{\frac{\log(1/\delta)}{n}} + \frac{2^{3/4}\eta K_g}{\sqrt{\lambda n}}$$

where $c_0, c_1, c_2 > 0$ are some absolute constants, λ is from the regularity condition in Assumption 3, and $K_g = 4B/\eta + 2F$ where B is from the assumption that the policy parameter set is bounded in Assumption 2, and F is from the Assumption 1.

Corollary 3 ("Fast" Estimation error of θ^{χ^2}). *Let $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$*

$$\|\hat{\theta}_n^{\chi^2} - \theta^{\chi^2}\|_2 \leq \frac{16c_0K_g\eta}{\lambda} \sqrt{\frac{d}{n}} + \left(\frac{8c_1K_g\eta}{\lambda} + \frac{2K_g\eta\sqrt{c_2}}{\sqrt{\lambda}} \right) \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\sqrt{2}\eta K_g}{\sqrt{\lambda n}}$$

where $c_0, c_1, c_2 > 0$ are some absolute constants, λ is from the regularity condition in Assumption 3, and $K_g = 4B/\eta + 2F$ where B is from the assumption that the policy parameter set is bounded in Assumption 2, and F is from the Assumption 1.

Our fast-rate analysis shows that DRO-REBEL achieves the minimax-optimal $O(n^{-1/2})$ estimation error, even under a Wasserstein, KL, or χ^2 ambiguity set, simply by combining strong convexity with a localized Rademacher complexity argument. This matches the classical parametric rate in M-estimation theory and parallels the finite-sample guarantees for generic Wasserstein DRO obtained by Gao [2022]. Crucially, the same localized-complexity machinery could be applied to other DRO settings under similar assumptions. In this way, DRO-REBEL both tightens constants in the slow regime and offers a clear recipe for restoring the optimal $n^{-1/2}$ convergence rate for common robust preference-learning algorithms.

4.4 Distributionally Robust DPO "Fast Rates"

Crucially, the same localized-complexity machinery could be applied to the WDPO and KLDPO analyses of Xu et al. [2025] to upgrade their $n^{-1/4}$ rates to $n^{-1/2}$ —albeit with larger constants, since REBEL's relative-reward regression loss is fundamentally simpler (no logistic link) and has smaller Lipschitz and curvature parameters. We significantly advance the theoretical understanding of robust DPO, specifically for Wasserstein DPO (WDPO) and KL DPO (KLDPO) and close a critical gap by establishing a superior rate of convergence compared to previous analyses, including that provided by Xu et al. [2025], with the following theorems:

Theorem 5 (Parametric $O(n^{-1/2})$ Rate for Wasserstein DPO). *Let the Wasserstein DPO empirical and population minimizers be defined as:*

$$\hat{\theta}_n^{\mathcal{W}_p} = \arg \min_{\theta \in \Theta} \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon_n), \quad \theta^{\mathcal{W}_p} = \arg \min_{\theta \in \Theta} \mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon_n)$$

where $\mathcal{L}^{\mathcal{W}_p}$ is the Wasserstein distributionally robust DPO objective with ambiguity radius ε_n . Assume all the assumptions made in the previous analyses. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the estimation error of the Wasserstein DPO estimator is bounded by:

$$\|\hat{\theta}_n^{\mathcal{W}_p} - \theta^{\mathcal{W}_p}\|_2 \leq \frac{32c_0\beta}{\gamma\lambda} \sqrt{\frac{d}{n}} + \frac{16c_1\beta}{\gamma\lambda} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8L_{\ell,z}}{\gamma\lambda n}} + \sqrt{\frac{4c_2 \log(1 + e^{4\beta B}) \log(1/\delta)}{\gamma\lambda n}}$$

where $\gamma = \frac{\beta^2 e^{4\beta B}}{(1 + e^{4\beta B})^2}$ is the DPO curvature constant, c_0, c_1, c_2 are universal constants from concentration inequalities, and $L_{\ell,z}$ is from Assumption 5.

Theorem 6 (Parametric $O(n^{-1/2})$ Rate for KL-DPO). *Let the KL-DPO empirical and population minimizers be defined as:*

$$\hat{\theta}_n^{\text{KL}} = \arg \min_{\theta \in \Theta} \mathcal{L}_n^{\text{KL}}(\theta; \varepsilon_n), \quad \theta^{\text{KL}} = \arg \min_{\theta \in \Theta} \mathcal{L}^{\text{KL}}(\theta; \varepsilon_n)$$

where \mathcal{L}^{KL} is the KL-divergence distributionally robust DPO objective with ambiguity radius ε_n . Assume all the assumptions made in the previous analyses. Then for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the estimation error of the KL-DPO estimator is bounded by:

$$\|\hat{\theta}_n^{\text{KL}} - \theta^{\text{KL}}\|_2 \leq \frac{32c_0\beta}{\gamma\lambda} \sqrt{\frac{d}{n}} + \frac{16c_1\beta}{\gamma\lambda} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8\log(1 + e^{4\beta B})}{\gamma\lambda n}} + \sqrt{\frac{4c_2 \log(1 + e^{4\beta B}) \log(1/\delta)}{\gamma\lambda n}}$$

where $\gamma = \frac{\beta^2 e^{4\beta B}}{(1 + e^{4\beta B})^2}$ and c_0, c_1, c_2 are universal constants.

We defer readers to the Appendix for the proofs, particularly Appendix J and K

5 Related Work and Comparisons

Our work on Distributionally Robust REBEL (DRO-REBEL) contributes to the growing body of literature on robust reinforcement learning from human feedback (RLHF), particularly concerning out-of-distribution (OOD) generalization and sample efficiency. We find strong theoretical alignment with recent advancements in χ^2 -based preference learning and offer distinct advantages in sample complexity compared to other distributionally robust RLHF approaches.

5.1 χ^2 -Based Preference Learning

Very recently, Huang et al. [2025] introduced χ^2 -PO, a one-line modification of Direct Preference Optimization (DPO) that replaces the usual log-link with a mixed $\chi^2 + \text{KL}$ link and implicitly enforces pessimism via the χ^2 -divergence. Their main result shows that χ^2 -PO achieves a sample-complexity guarantee scaling as

$$J(\pi^*) - J(\hat{\pi}) \lesssim \sqrt{\frac{C_{\pi^*} \log(|\Pi|/\delta)}{n}}$$

where $C_{\pi^*} = 1 + 2D_{\chi^2}(\pi^* \|\pi_{\text{ref}})$ is the single-policy concentrability coefficient (cf. Theorem 3.1 of Huang et al. [2025]). Our "fast rate" analysis for the χ^2 -robust REBEL update (Corollary 3) recovers exactly the same $O(n^{-1/2})$ parametric rate under analogous linear-policy, data-coverage, and strong-convexity assumptions. In both cases the key is that χ^2 -regularization induces a heavy-tailed density-ratio barrier and uniform quadratic growth, allowing a localized Rademacher complexity argument to restore the minimax $n^{-1/2}$ rate. Thus, the theoretical insights of Huang et al. [2025] on the power of χ^2 -divergence to suppress overoptimization are fully consistent with—and indeed validated by—our fast-rate guarantees for χ^2 -REBEL as a sample-efficient and stable optimizer under preference noise.

5.2 Comparison with Distributionally Robust RLHF (Mandal et al., 2025)

Our work on DRO-REBEL directly addresses the robust policy optimization problem by minimizing a distributionally robust REBEL objective, assuming a fixed reward model. A parallel effort is made by Mandal et al. [2025], who also tackle robust policy optimization within the RLHF context. A key distinction lies in the choice of ambiguity sets and the resulting sample complexity guarantees for the robust policy. Mandal et al. [2025] investigate two main algorithms for robust policy optimization, both primarily using Total Variation (TV) distance for their ambiguity sets:

1. **Robust Natural Policy Gradient (NPG):** Their Theorem 2 states that to achieve an ε -optimal robust value function, their robust NPG algorithm requires $T = O(\log(1/\varepsilon))$ iterations. However, each iteration demands a large minibatch size of $n \geq O(\log(\log \log(1/\varepsilon))/\varepsilon^4)$. This translates to a total sample complexity of $O(\log(1/\varepsilon) \cdot \log(\log \log(1/\varepsilon))/\varepsilon^4)$, indicating a sub-polynomial dependence on $1/\varepsilon$.
2. **Robust Direct Preference Optimization (DPO):** For their robust DPO algorithm, Theorem 3 from Mandal et al. [2025] states that $T = O(1/\varepsilon^2)$ iterations are needed, with each iteration requiring a minibatch size n such that $n/\log n \geq O(1/\varepsilon^2)$. This results in a total sample complexity of $O(\log(1/\varepsilon)/\varepsilon^4)$ for achieving an ε -optimal robust DPO loss.

In contrast, our DRO-REBEL framework provides “fast rate” sample complexity guarantees for the parameter estimation error $\|\hat{\theta}_n - \theta^*\|_2$ across various ambiguity sets (Corollaries 1, 2, and 3). For all considered divergences (Wasserstein,

KL, and χ^2), and under appropriate choices of the ball radius ε_n (e.g., $\varepsilon_n = O(n^{-1})$ for Wasserstein and $O(n^{-2})$ for KL and χ^2), our results show that the estimation error converges at a parametric $O(n^{-1/2})$ rate. This implies that to achieve an estimation error of $O(\varepsilon)$, our method requires a total sample size of $n = O(1/\varepsilon^2)$. This also applies to our analysis of robust DPO (Theorem 5 and 6)

Our derived “fast rates” are significantly more efficient in terms of sample complexity compared to the results presented by Mandal et al. [2025] for their TV-distance based robust policy optimization algorithms. The $O(1/\varepsilon^2)$ sample complexity to achieve an $O(\varepsilon)$ parameter estimation error (equivalent to an $O(n^{-1/2})$ rate) aligns with the minimax optimal rates often observed in parametric statistical problems. Our theoretical framework demonstrates that certain divergence choices in DRO-REBEL and simplifying the structure and approach of the problem to relative regression can lead to significantly improved sample complexities for robust policy learning.

Algorithm / Work	Parameter	Ambiguity Set(s)	Error Rate (n)	Sample Complexity (ε)
Our Robust DPO	θ	\mathcal{W}_p, KL	$O(n^{-1/2})$	$O(1/\varepsilon^2)$
Mandal et al. (Robust DPO)	Policy	TV	$O(n^{-1/4})$	$O(\log(1/\varepsilon)/\varepsilon^4)$
Xu et al. (Robust DPO)	θ	\mathcal{W}_p, KL	$O(n^{-1/4})$	$O(1/\varepsilon^4)$

Table 1: Comparison of theoretical convergence rates and sample complexities for robust DPO algorithms.

6 Approximate Tractable Algorithms for Robust LLM Alignment

While our Distributionally Robust REBEL (DRO-REBEL) formulations benefit from finite-sample guarantees which are minimax optimal, directly solving the minimax objective using stochastic gradient descent methods can be computationally challenging. As Xu et al. [2025] also point out in the context of robust DPO, this challenge arises because we do not have direct control over the data distribution $\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; D)$ within the uncertainty set, as it is not parameterized in a straightforward manner. Furthermore, the preference data are generated according to the nominal distribution \mathbb{P}° , meaning we lack samples from any other distributions within the uncertainty set $\mathcal{B}_\varepsilon(\mathbb{P}^\circ; D)$. To overcome this, we introduce principled tractable algorithms that approximate the solution to our DRO-REBEL objectives. Our algorithms for solving Wasserstein-DRO-REBEL and KL-DRO-REBEL are largely the same as those of Xu et al. [2025]’s WDPO and KLDPO. However, we will derive and propose an algorithm for χ^2 -DRO-REBEL that can efficiently be solved using stochastic gradient descent methods.

6.1 Tractable Wasserstein DRO-REBEL (WD-REBEL)

The connection between Wasserstein distributionally robust optimization (DRO) and regularization has been established previously in the literature, see Shafieezadeh-Abadeh et al. [2019] for example. We leverage recent progress in Wasserstein theory on connecting Wasserstein DRO to regularization. For p -Wasserstein DRO, $p \in (1, \infty]$, Gao and Kleywegt [2022] (Theorem 1) shows that for a broad class of loss functions (potentially non-convex and non-smooth), with high probability, Wasserstein DRO is asymptotically equivalent to a variation regularization. In particular, an immediate consequence is that, when $p = 2$:

$$\min_{\theta \in \Theta} \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}_n^\circ; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] = \min_{\theta \in \Theta} \left\{ \mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell(z; \theta)] + \varepsilon_n \sqrt{(1/n) \sum_{i=1}^n \|\nabla_z \ell(z_i; \theta)\|_2^2} \right\} + O_p(1/n),$$

where $\rho_n = O(1/\sqrt{n})$. This indicates that one can approximately solve the Wasserstein DRO objective by adding a gradient regularization term to the empirical risk minimization (ERM) loss, $\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell(z; \theta)]$. Based on this, we propose a tractable WD-REBEL algorithm in Algorithm 2.

6.2 Tractable KL-DRO-REBEL (KL-REBEL)

We utilize the following proposition established by Xu et al. [2025] to show that we can approximate the worst-case probability distribution in a KL uncertainty set with respect to a given loss function.

Proposition 1 (Worst-case distribution). *Let $\mathbb{P} \in \mathbb{R}^n$ be the worst-case distribution with respect to a loss function ℓ and KL uncertainty around the empirical distribution \mathbb{P}_n° , defined as $\mathbb{P} = \sup_{\mathbb{P}: D_{\text{KL}}(\mathbb{P} \parallel \mathbb{P}_n^\circ) \leq \rho} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$. The worst-case*

Algorithm 2 WD-REBEL Algorithm

Require: Dataset $\mathcal{D} = \{z_i\}_{i=1}^n$ (e.g., preference pairs), reference policy π_{ref} , robustness hyperparameter ρ_0 , learning rate η , initial policy π_θ .

- 1: **while** θ has not converged **do**
 - 2: Calculate the non-robust REBEL loss $\ell(z_i; \theta)$ for each $z_i \in \mathcal{D}$.
 - 3: Calculate the non-robust empirical REBEL loss $L_{\text{REBEL}}(\pi_\theta; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \ell(z_i; \theta)$.
 - 4: Calculate the gradient regularizer term: $R(\pi_\theta; \mathcal{D}) = \rho_0 \left(\frac{1}{n} \sum_{i=1}^n \|\nabla_{z_i} \ell(z_i; \theta)\|_2^2 \right)^{1/2}$.
 - 5: Calculate the approximate WD-REBEL loss: $L_W(\theta, \rho_0) = L_{\text{REBEL}}(\pi_\theta; \mathcal{D}) + R(\pi_\theta; \mathcal{D})$.
 - 6: $\theta \leftarrow \theta - \eta \nabla_\theta L_W(\theta, \rho_0)$.
 - 7: **end while**
 - 8: **return** π_θ .
-

distribution \mathbb{P} is related to \mathbb{P}_n° through

$$\mathbb{P}(i) \propto \mathbb{P}_n(i) \cdot \exp \left(\frac{1}{\tau} (\ell(z_i; \theta) - \sum_{j=1}^n \mathbb{P}_n(j) \ell(z_j; \theta)) \right),$$

where $\tau > 0$ is some constant.

A proof of this proposition can be found in Appendix D of Xu et al. [2025]. Based on Proposition 1, we propose a tractable KL-REBEL algorithm in Algorithm 3.

Algorithm 3 KL-REBEL Algorithm

Require: Dataset $\mathcal{D} = \{z_i\}_{i=1}^n$, reference policy π_{ref} , robustness temperature parameter τ , learning rate η , initial policy π_θ .

- 1: **while** θ has not converged **do**
 - 2: Calculate the non-robust REBEL loss $\ell(z_i; \theta)$ for each $z_i \in \mathcal{D}$.
 - 3: Approximate the worst-case weights assuming $\mathbb{P}_n(i) = 1/n$: $\tilde{\mathbb{P}}(i) = \exp \left(\frac{1}{\tau} \left(\ell(z_i; \theta) - \frac{1}{n} \sum_{j=1}^n \ell(z_j; \theta) \right) \right)$.
 - 4: Normalize the weights: $\mathbb{P}(i) = \frac{\tilde{\mathbb{P}}(i)}{\sum_{k=1}^n \tilde{\mathbb{P}}(k)}$.
 - 5: Calculate the approximate KL-REBEL loss: $L_{\text{KL}}(\theta; \mathcal{D}) = \sum_{i=1}^n \mathbb{P}(i) \cdot \ell(z_i; \theta)$.
 - 6: $\theta \leftarrow \theta - \eta \nabla_\theta L_{\text{KL}}(\theta, \rho)$.
 - 7: **end while**
 - 8: **return** π_θ .
-

6.3 Tractable χ^2 -DRO-REBEL (χ^2 -REBEL)

We exploit the dual formulation of the χ^2 -DRO objective (e.g. Namkoong and Duchi [2017a]) to obtain a one-dimensional inner solve and closed-form worst-case weights.

Proposition 2 (Dual form & worst-case weights). *Let $\ell_i = \ell(z_i; \theta)$ for $i = 1, \dots, n$ and set $\mathcal{L}_n^{\chi^2}(\theta; \varepsilon_n) = \sup_{\mathbb{P}: D_{\chi^2}(\mathbb{P} \parallel \mathbb{P}_n^\circ) \leq \varepsilon_n} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)]$. Then*

$$\mathcal{L}_n^{\chi^2}(\theta; \varepsilon_n) = \inf_{\eta \in \mathbb{R}} \left\{ \eta + \sqrt{\frac{2\varepsilon_n}{n} \sum_{i=1}^n (\ell_i - \eta)_+^2} \right\},$$

and the supremum is attained by the “re-weighted” empirical distribution

$$w_i = \frac{(\ell_i - \eta^*)_+}{n \lambda^*}, \quad \lambda^* = \sqrt{\frac{2\varepsilon_n}{n} \sum_{j=1}^n (\ell_j - \eta^*)_+^2},$$

where η^* minimizes the inner objective.

We defer the proof to Appendix L. Based on Proposition 2, we propose a tractable χ^2 -REBEL Algorithm in Algorithm 4

Algorithm 4 χ^2 -REBEL Algorithm

Require: Dataset $\mathcal{D} = \{z_i\}_{i=1}^n$, robustness radius ρ , learning rate α , initial policy π_θ

- 1: **while** θ not converged **do**
- 2: Compute pointwise REBEL loss $\ell_i \leftarrow \ell(z_i; \theta)$ and gradients $g_i \leftarrow \nabla_\theta \ell(z_i; \theta)$
- 3: (*Inner 1-D solve*) find

$$\eta^* = \arg \min_{\eta \in \mathbb{R}} \left\{ \eta + \sqrt{\frac{2\rho}{n} \sum_{i=1}^n (\ell_i - \eta)_+^2} \right\}$$

via sorting $\{\ell_i\}$ and binary search

- 4: Compute $\lambda^* = \sqrt{\frac{2\rho}{n} \sum_{i=1}^n (\ell_i - \eta^*)_+^2}$
 - 5: Form weights $w_i = \frac{(\ell_i - \eta^*)_+}{n \lambda^*}$
 - 6: Robust gradient: $G = \sum_{i=1}^n w_i g_i$
 - 7: Gradient step: $\theta \leftarrow \theta - \alpha G$
 - 8: **end while**
 - 9: **return** π_θ
-

Each outer step of Algorithm 4 has a computational cost of $O(n \log n + nd)$. This complexity arises from the two primary components of the update: the inner one-dimensional solve for the optimal dual variable η^* and the subsequent robust gradient computation and parameter update. While the inner solve is an optimization over the unbounded domain $\eta \in \mathbb{R}$, its specific structure allows for a highly efficient solution. Let the inner objective function be denoted by $f(\eta)$:

$$f(\eta) = \eta + \sqrt{\frac{2\rho}{n} \sum_{i=1}^n (\ell_i - \eta)_+^2}$$

where $(\cdot)_+ := \max(0, \cdot)$. This function $f(\eta)$ can be shown to be **convex** in η . It is the sum of a linear function, η , and a second term which is a composition of the convex L2-norm with a vector of convex functions, $v_i(\eta) = (\ell_i - \eta)_+$. The convexity of $f(\eta)$ ensures the existence of a unique global minimum and implies that its subgradient, $\partial f(\eta)$, is monotonically non-decreasing.

This monotonicity allows for the use of an efficient binary search to find the unique η^* where $0 \in \partial f(\eta^*)$. The search domain is not the entirety of \mathbb{R} ; one can show that the optimal solution is bounded, i.e., $\eta^* \in (-\infty, \max_i \ell_i]$, because for any $\eta > \max_i \ell_i$, the objective simplifies to $f(\eta) = \eta$, which is strictly increasing. The cost of this inner step is therefore dominated by the initial sorting of the loss values $\{\ell_i\}_{i=1}^n$ to identify the function’s non-differentiable “breakpoints,” which takes $O(n \log n)$ time.

With the inner problem solved tractably, we consider the outer optimization over θ . Since the pointwise loss $\ell(z; \theta)$ is uniformly bounded and Lipschitz continuous (Appendix B), and the overall robust objective $\mathcal{L}_n^{\chi^2}(\theta; \rho)$ is strongly convex (Appendix E), we can guarantee that standard optimization algorithms (e.g., projected gradient descent) will converge to the unique optimal policy parameters θ^* . This guarantee of convergence also holds for the Wasserstein and KL-divergence counterparts of our algorithm.

7 Conclusion

In this work, we introduced *DRO-REBEL*, a unified family of distributionally-robust variants of the REBEL framework for offline RLHF. By instantiating ambiguity sets via KL, χ^2 , and type- p Wasserstein divergences and exploiting strong/Fenchel duality, each robust policy update reduces to a simple relative-reward regression. Under standard linear-policy and data-coverage assumptions, we proved that all DRO-REBEL variants achieve an $O(n^{-1/4})$ “slow” estimation-error rate with substantially tighter constants than prior DRO-DPO methods and, via a localized Rademacher-complexity argument, recover the minimax-optimal $O(n^{-1/2})$ “fast” parametric rate. Our analysis shows that DRO-REBEL not only tightens slow-rate bounds but also restores the classical parametric rate under distributional shifts.

Along the way, we learned that strong/Fenchel duality can collapse complex DRO updates into tractable regressions, that improving constant factors in slow-rate bounds can have outsized practical impact and, perhaps most surprisingly,

that worst-case robustness and optimal $O(n^{-1/2})$ convergence need not be at odds. This work also raises several open questions: How should practitioners choose or adapt between KL, χ^2 , and Wasserstein ambiguity sets in real RLHF pipelines? Can the same dual-regression approach and fast-rate analysis be extended to neural (nonlinear) policy classes? And how might one jointly integrate robust reward modeling with DRO-REBEL to further bolster alignment under noisy or adversarial feedback?

Looking ahead, we plan to validate these theoretical findings empirically. We will train our robust RLHF algorithms on the Unified-Feedback dataset [Jiang et al., 2023] and evaluate out-of-distribution robustness on established reward-evaluation benchmarks such as Reward-Bench [Lambert et al., 2024], MT-Bench [Zheng et al., 2023], and HHH-Alignment [Askell et al., 2021]. With more time, we would also conduct systematic hyperparameter sweeps (e.g. on ε , η , dual bounds), prototype DRO-REBEL with transformer-parameterized policies, and explore adaptive schemes for selecting or mixing ambiguity sets online. We anticipate that these experiments will corroborate our theoretical guarantees, demonstrating superior generalization under preference shifts and resistance to overoptimization, thereby providing a unified, principled methodology for reliably aligning LLMs with diverse and uncertain human preferences. Due to time constraints, we were unable to complete these experiments before submission. However, we are actively working on them over the summer and expect the empirical results to align closely with our theory.

References

- Alekh Agarwal, Nan Jiang, Sham M. Kakade, and Wen Sun. *Reinforcement Learning: Theory and Algorithms*. 2021a. URL <https://rltheorybook.github.io/>.
- Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021b. URL <http://jmlr.org/papers/v22/19-736.html>.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. A general language assistant as a laboratory for alignment, 2021. URL <https://arxiv.org/abs/2112.00861>.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislaw Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022. URL <https://arxiv.org/abs/2204.05862>.
- A. Beck. *First-Order Methods in Optimization*. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics, 2017. ISBN 9781611974997. URL <https://books.google.com/books?id=wrk4DwAAQBAJ>.
- A. Beck. *Introduction to Nonlinear Optimization: Theory, Algorithms and Applications with Python and MATLAB*. MOS-SIAM series on optimization. Society for Industrial and Applied Mathematics, 2023. ISBN 9781611977615. URL <https://books.google.com/books?id=YDrizwEACAAJ>.
- S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford, 2013. ISBN 9780199535255. URL <https://books.google.com/books?id=5oo4YIz6tR0C>.
- Olivier Bousquet. A bennett concentration inequality and its application to suprema of empirical processes. *Comptes Rendus Mathématique*, 334(6):495–500, 2002. ISSN 1631-073X. doi: [https://doi.org/10.1016/S1631-073X\(02\)02292-6](https://doi.org/10.1016/S1631-073X(02)02292-6). URL <https://www.sciencedirect.com/science/article/pii/S1631073X02022926>.
- Alexander Bukharin, Ilgee Hong, Haoming Jiang, Zichong Li, Qingru Zhang, Zixuan Zhang, and Tuo Zhao. Robust reinforcement learning from corrupted human feedback, 2024. URL <https://arxiv.org/abs/2406.15568>.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Biyik, Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. Open problems and fundamental limitations of reinforcement learning from human feedback, 2023. URL <https://arxiv.org/abs/2307.15217>.

- Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Amrit Singh Bedi, and Mengdi Wang. Maxmin-rlhf: Alignment with diverse human preferences, 2024. URL <https://arxiv.org/abs/2402.08925>.
- Sayak Ray Chowdhury, Anush Kini, and Nagarajan Natarajan. Provably robust dpo: Aligning language models with noisy feedback, 2024. URL <https://arxiv.org/abs/2403.00409>.
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences, 2023. URL <https://arxiv.org/abs/1706.03741>.
- John Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization, 2020. URL <https://arxiv.org/abs/1810.08750>.
- John C. Duchi and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. *Mathematics of Operations Research*, 47(2):753–789, 2022.
- Rui Gao. Finite-sample guarantees for wasserstein distributionally robust optimization: Breaking the curse of dimensionality, 2022. URL <https://arxiv.org/abs/2009.04382>.
- Rui Gao and Anton J. Kleywegt. Distributionally robust stochastic optimization with wasserstein distance, 2022. URL <https://arxiv.org/abs/1604.02199>.
- Yang Gao, Yuxin Xie, Nan Jiang, and Lihong Wang. Distributionally robust policy evaluation and learning in offline contextual bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 8512–8530, 2022.
- Zhaolin Gao, Jonathan D. Chang, Wenhao Zhan, Owen Oertell, Gokul Swamy, Kianté Brantley, Thorsten Joachims, J. Andrew Bagnell, Jason D. Lee, and Wen Sun. Rebel: Reinforcement learning via regressing relative rewards, 2024. URL <https://arxiv.org/abs/2404.16767>.
- Ilgee Hong, Zichong Li, Alexander Bukharin, Yixiao Li, Haoming Jiang, Tianbao Yang, and Tuo Zhao. Adaptive preference scaling for reinforcement learning with human feedback, 2024. URL <https://arxiv.org/abs/2406.02764>.
- Audrey Huang, Wenhao Zhan, Tengyang Xie, Jason D. Lee, Wen Sun, Akshay Krishnamurthy, and Dylan J. Foster. Correcting the mythos of kl-regularization: Direct alignment without overoptimization via chi-squared preference optimization, 2025. URL <https://arxiv.org/abs/2407.13399>.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion, 2023. URL <https://arxiv.org/abs/2306.02561>.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl?, 2022. URL <https://arxiv.org/abs/2012.15085>.
- Sham M Kakade. A natural policy gradient. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001. URL https://proceedings.neurips.cc/paper_files/paper/2001/file/4b86abe48d358ecf194c56c69108433e-Paper.pdf.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of rlhf on llm generalisation and diversity, 2024. URL <https://arxiv.org/abs/2310.06452>.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Rewardbench: Evaluating reward models for language modeling, 2024. URL <https://arxiv.org/abs/2403.13787>.
- E. L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer, 2nd edition, 1998. Chapter 5 discusses Barankin-type bounds, including HCR.
- Will LeVine, Benjamin Pikus, Anthony Chen, and Sean Hendryx. A baseline analysis of reward models’ ability to accurately analyze foundation models under distribution shift, 2024. URL <https://arxiv.org/abs/2311.14743>.
- Debmalya Mandal, Paulius Sasnauskas, and Goran Radanovic. Distributionally robust reinforcement learning with human feedback, 2025. URL <https://arxiv.org/abs/2503.00539>.

- Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 2010–2020. PMLR, 26–28 Aug 2020. URL <https://proceedings.mlr.press/v108/modi20a.html>.
- Hongseok Namkoong and John C Duchi. Variance-based regularization with convex objectives. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017a. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/5a142a55461d5fef016acfb927fee0bd-Paper.pdf.
- Hongseok Namkoong and John C. Duchi. Variance-based regularization with convex objectives. In *Advances in Neural Information Processing Systems*, volume 30, 2017b.
- Andi Nika, Debmalya Mandal, Parameswaran Kamalaruban, Georgios Tzannetos, Goran Radanović, and Adish Singla. Reward model learning vs. direct policy optimization: A comparative analysis of learning from human preferences, 2024. URL <https://arxiv.org/abs/2403.01857>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- Vishakh Padmakumar, Chuanyang Jin, Hannah Rose Kirk, and He He. Beyond the binary: Capturing diverse preferences with reward regularization, 2024. URL <https://arxiv.org/abs/2412.03822>.
- Yury Polyanskiy. Lecture notes on information theory, chapter 29, ece563 (uiuc). Technical report, Department of Electrical and Computer Engineering, University of Illinois at Urbana–Champaign, 2017. URL https://people.lids.mit.edu/yp/homepage/data/LN_stats.pdf. Archived (PDF) from the original on 2022-05-24. Retrieved 2022-05-24.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- Soroosh Shafieezadeh-Abadeh, Daniel Kuhn, and Peyman Mohajerin Esfahani. Regularization via mass transportation, 2019. URL <https://arxiv.org/abs/1710.10016>.
- Harvineet Singh Shah, Michael Jung, Kyomin Jung, and Hwanjo Kim. Robust optimization for fairness with noisy protected groups. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5702–5709, 2020.
- Alexander Shapiro and Huan Xu. Wasserstein distributionally robust optimization: Theory and applications in machine learning. *Operations Research*, 70(5):3007–3030, 2022.
- A. van der Vaart and J.A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, 1996. ISBN 9780387946405. URL <https://books.google.com/books?id=0CenCW9qmp4C>.
- Ramon van Handel. Probability in high dimensions. 2016. URL <https://web.math.princeton.edu/~rvan/APC550.pdf>.
- Ruosong Wang, Dean P. Foster, and Sham M. Kakade. What are the statistical limits of offline rl with linear function approximation?, 2020. URL <https://arxiv.org/abs/2010.11895>.
- Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jiawei Chen, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. Towards robust alignment of language models: Distributionally robustifying direct preference optimization, 2024. URL <https://arxiv.org/abs/2407.07880>.
- Zaiyan Xu, Sushil Vemuri, Kishan Panaganti, Dileep Kalathil, Rahul Jain, and Deepak Ramachandran. Distributionally robust direct preference optimization, 2025. URL <https://arxiv.org/abs/2502.01930>.

- Yuzi Yan, Xingzhou Lou, Jialian Li, Yiping Zhang, Jian Xie, Chao Yu, Yu Wang, Dong Yan, and Yuan Shen. Reward-robust rlhf in llms, 2024. URL <https://arxiv.org/abs/2409.15360>.
- Fan Yang, Shixiang Gu, Sergey Levine, Dale Schuurmans, and Ofir Nachum. Wasserstein distributional reinforcement learning. In *International Conference on Learning Representations*, 2021.
- Xiaoying Zhang, Jean-Francois Ton, Wei Shen, Hongning Wang, and Yang Liu. Overcoming reward overoptimization via adversarial policy optimization with lightweight uncertainty estimation, 2024. URL <https://arxiv.org/abs/2403.05171>.
- Yichen Zhang, Yuxin Xie, and Lihong Wang. Generalization in reinforcement learning with stochastic mirror descent. In *International Conference on Machine Learning*, volume 119, pages 11188–11197, 2020.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- Yingxue Zhou, Haoran Li, Longbo Huang, and Peilin Zhao. Distributionally robust exploration in multi-armed bandits. In *Advances in Neural Information Processing Systems*, volume 35, pages 17328–17340, 2022.
- Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 43037–43067. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/zhu23f.html>.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020. URL <https://arxiv.org/abs/1909.08593>.

A Auxiliary Technical Tools

A.1 Wasserstein Theory

Lemma 2 (Gao and Kleywegt [2022], Theorem 1; Strong Duality for DRO with Wasserstein Distance). *Consider any $p \in [1, \infty)$, any $\nu \in \mathcal{P}(\Xi)$, any $\rho > 0$, and any $\Psi \in L^1(\nu)$ such that the growth rate κ of Ψ satisfies*

$$\kappa := \inf \left\{ \eta \geq 0 : \int_{\Xi} \Phi(\eta, \zeta) \nu(d\zeta) > -\infty \right\} < \infty, \quad (13)$$

where

$$\Phi(\eta, \zeta) := \inf_{\xi \in \Xi} \{ \eta d^p(\xi, \zeta) - \Psi(\xi) \}.$$

Then strong duality holds with finite optimal value $v_p = v_D \leq \infty$, where the primal and dual problems are

$$v_p = \sup_{\mu \in \mathcal{P}(\Xi)} \left\{ \int_{\Xi} \Psi(\xi) \mu(d\xi) : W_p(\mu, \nu) \leq \rho \right\}, \quad (\text{Primal}) \quad (3)$$

$$v_D = \inf_{\eta \geq 0} \left\{ \eta \rho^p - \int_{\Xi} \inf_{\xi \in \Xi} [\eta d^p(\xi, \zeta) - \Psi(\xi)] \nu(d\zeta) \right\}. \quad (\text{Dual}) \quad (4)$$

Lemma 3 (Gao and Kleywegt [2022], Lemma 2(ii); Properties of the growth κ). *Suppose that $\nu \in \mathcal{P}_p(\Xi)$. Then the growth rate κ in (13) is finite if and only if there exists $\zeta_0 \in \Xi$ and constants $L, M > 0$ such that*

$$\Psi(\xi) - \Psi(\zeta_0) \leq L d^p(\xi, \zeta_0) + M, \quad \forall \xi \in \Xi. \quad (14)$$

Corollary 4. *Consider any bounded loss function ℓ over a bounded space Ξ . Then the duality in Lemma 2 holds.*

Proof. Immediate from Lemma 3 by choosing $L = \text{diam}(\Xi)^p$ and $M = \sup_{\xi \in \Xi} |\Psi(\xi)|$. \square

A.2 Optimization

Lemma 4 (Beck [2023], Theorem 1.24; Linear Approximation Theorem). *Let $f: U \rightarrow \mathbb{R}$ be twice continuously differentiable on an open set $U \subseteq \mathbb{R}^n$, and let $x, y \in U$ satisfy $[x, y] \subset U$. Then there exists $\xi \in [x, y]$ such that*

$$f(y) = f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2} (y - x)^\top \nabla^2 f(\xi) (y - x).$$

Lemma 5 (Beck [2017], Theorem 5.24; First-order characterizations of strong convexity). *Let $f: E \rightarrow (-\infty, \infty]$ be a proper, closed, convex function, and let $\sigma > 0$. The following are equivalent:*

1. *For all $x, y \in \text{dom}(f)$ and $\lambda \in [0, 1]$,*

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\sigma}{2} \lambda(1 - \lambda) \|x - y\|^2.$$

2. *For all $x \in \text{dom}(\partial f)$, $y \in \text{dom}(f)$ and $g \in \partial f(x)$,*

$$f(y) \geq f(x) + \langle g, y - x \rangle + \frac{\sigma}{2} \|y - x\|^2.$$

Lemma 6 (Beck [2017], Theorem 5.25; Existence and uniqueness of minimizer). *Let $f: E \rightarrow (-\infty, \infty]$ be proper, closed, and σ -strongly convex with $\sigma > 0$. Then:*

1. *f has a unique minimizer x^* .*

2. *For all $x \in \text{dom}(f)$,*

$$f(x) - f(x^*) \geq \frac{\sigma}{2} \|x - x^*\|^2.$$

A.3 Distributionally Robust Optimization

The f -divergence between distributions \mathbb{P} and \mathbb{P}_0 on \mathcal{X} is

$$D_f(P \| P_0) = \int_{\mathcal{X}} f\left(\frac{d\mathbb{P}}{d\mathbb{P}_0}\right) d\mathbb{P}_0, \quad (15)$$

where f is a convex function (e.g. $f(t) = t \log t$ gives KL divergence). For a loss $\ell: \mathcal{X} \rightarrow \mathbb{R}$:

Lemma 7 (Duchi and Namkoong [2020], Proposition 1). *Let D_f be as in (15). Then*

$$\sup_{P: D_f(P \| P_0) \leq \rho} \mathbb{E}_P[\ell(X)] = \inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \left\{ \lambda f^*\left(\frac{\ell(X) - \eta}{\lambda}\right) + \lambda \rho + \eta \right\}, \quad (16)$$

where $f^*(s) = \sup_{t \geq 0} \{st - f(t)\}$ is the Fenchel conjugate of f .

A.4 Empirical Process Theory

Lemma 8 (van der Vaart and Wellner [1996], Lemma 2.3.1; Symmetrization). *For every nondecreasing, convex $\Phi: \mathbb{R} \rightarrow \mathbb{R}$ and class of measurable functions \mathcal{F} ,*

$$\mathbb{E}^*[\Phi(\|P_n - P\|_{\mathcal{F}})] \leq \mathbb{E}^*[\Phi(2\|P_n^\circ\|_{\mathcal{F}})],$$

where the outer expectations \mathbb{E}^* are taken over the data generating distribution and Rademacher random variables and P_n° is the symmetrized process.

Corollary 5 (Boucheron et al. [2013], Corollary 13.2; Dudley's Entropy Integral). *Let \mathcal{T} be a finite pseudometric space and let $(X_t)_{t \in \mathcal{T}}$ be a collection of random variables such that, for all $t, t' \in \mathcal{T}$ and all $\lambda > 0$,*

$$\log \mathbb{E}[e^{\lambda(X_t - X_{t'})}] \leq \frac{\lambda^2 d^2(t, t')}{2}.$$

Then for any fixed $t_0 \in \mathcal{T}$, if we set

$$\delta = \sup_{t \in \mathcal{T}} d(t, t_0) \quad \text{and} \quad H(u, \mathcal{T}) = \log N(u, \mathcal{T}, d)$$

denoting the covering-number entropy at scale u , it holds that

$$\mathbb{E}\left[\sup_{t \in \mathcal{T}} (X_t - X_{t_0})\right] \leq 12 \int_0^{\delta/2} \sqrt{H(u, \mathcal{T})} \, du.$$

Theorem 7 (Bousquet [2002], Theorem 2.1). *Let $c > 0$, let X_i be independent random variables with distribution P , and let \mathcal{F} be a class of functions $f : X \rightarrow \mathbb{R}$. Assume that for all $f \in \mathcal{F}$,*

$$\mathbb{E}[f(X_i)] = 0, \quad \|f\|_\infty \leq c.$$

Let $\sigma > 0$ satisfy

$$\sigma^2 \geq \sup_{f \in \mathcal{F}} \text{Var}(f(X_i)).$$

Define

$$Z = \sup_{f \in \mathcal{F}} \sum_{i=1}^n f(X_i), \quad v = n\sigma^2 + 2c \mathbb{E}[Z], \quad h(u) = (1+u) \ln(1+u) - u.$$

Then for any $x \geq 0$,

$$\Pr(Z \geq \mathbb{E}[Z] + x) \leq \exp\left(-v h\left(\frac{x}{c v}\right)\right).$$

Moreover, with probability at least $1 - e^{-x}$,

$$Z \leq \mathbb{E}[Z] + \sqrt{2xv} + \frac{cx}{3}.$$

A.5 Variational Inequalities

Lemma 9 (van Handel [2016], Lemma 4.10; Gibbs Variational Principle). *Let $\mu, \nu \in \mathcal{P}(\Xi)$ be Borel probability measures supported on Ξ . Then*

$$\log \mathbb{E}_\mu [e^f] = \sup_\nu \{ \mathbb{E}_\nu [f] - D_{\text{KL}}(\mu \parallel \nu) \}$$

Theorem 8 ([Polyanskiy, 2017, Lehmann and Casella, 1998]; Hammersley-Chapman-Robbins (HCR) lower bound). *Let Θ be the set of parameters for a family of probability distributions $\{\mu_\theta : \theta \in \Theta\}$ on a sample space Ω . For any $\theta, \theta' \in \Theta$, let $\chi^2(\mu_{\theta'}; \mu_\theta)$ denote the χ^2 -divergence from μ_θ to $\mu_{\theta'}$. For any scalar random variable $\hat{g} : \Omega \rightarrow \mathbb{R}$ and any $\theta, \theta' \in \Theta$, we have*

$$\text{Var}_\theta[\hat{g}] \geq \sup_{\substack{\theta' \neq \theta \\ \theta' \in \Theta}} \frac{(\mathbb{E}_{\theta'}[\hat{g}] - \mathbb{E}_\theta[\hat{g}])^2}{\chi^2(\mu_{\theta'}; \mu_\theta)}.$$

A.6 Concentration Inequalities

Lemma 10 (Boucheron et al. [2013], Lemma 2.2; Hoeffding's Lemma). *Let Y be a random variable with $\mathbb{E}[Y] = 0$ and almost surely $Y \in [a, b]$. Define $\psi_Y(\lambda) = \log \mathbb{E}[e^{\lambda Y}]$. Then for all $\lambda \in \mathbb{R}$, $\psi_Y''(\lambda) \leq \frac{(b-a)^2}{4}$ and consequently Y is sub-Gaussian with proxy variance $(b-a)^2/4$, i.e. $Y \sim \mathcal{SG}\left(\frac{b-a}{2}\right)$.*

Using Hoeffding's lemma, one can prove Hoeffding's inequality using a standard Chernoff bound argument.

Lemma 11 (Hoeffding's inequality). *Let X_1, \dots, X_n be independent with $X_i \in [a_i, b_i]$ almost surely, and define*

$$S = \sum_{i=1}^n (X_i - \mathbb{E}[X_i]).$$

Then for every $t > 0$,

$$\mathbb{P}(S \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

In particular, if X_1, \dots, X_n are i.i.d. with mean μ and support $[a, b]$, then for all $t > 0$

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X]\right| \geq t\right) \leq 2 \exp\left(-\frac{2nt^2}{(b-a)^2}\right).$$

B Uniform and Lipschitzness bounds of $\ell(z; \theta)$

We first prove that $\ell(z; \theta)$ is uniformly bounded

Lemma 12 (Uniform bound on $\ell(z; \theta)$). *Let $K_g = \sup_{z, \theta} |g(z; \theta)| \leq 4B/\eta + 2F$ where $\ell(z; \theta) = g(z; \theta)^2$ with $z = (x, a^1, a^2) \sim \mathbb{P}^\circ$. Then $\sup_{z, \theta} |\ell(z; \theta)| = K_\ell = K_g^2$*

Proof of Lemma 12. Since we have that $\pi_\theta, \pi_{\theta_t} \in \Pi$, notice that

$$\begin{aligned} \log \left(\frac{\pi_\theta(a | x)}{\pi_{\theta_t}(a | x)} \right) &= \log \pi_\theta(a | x) - \log \pi_{\theta_t}(a | x) \\ &= \log \left(\frac{\exp(\theta^\top \psi(x, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta^\top \psi(x, a'))} \right) - \log \left(\frac{\exp(\theta_t^\top \psi(x, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta_t^\top \psi(x, a'))} \right) \\ &= \log \left(\exp((\theta - \theta_t)^\top \psi(x, a)) \right) + \log \left(\frac{\sum_{a' \in \mathcal{A}} \exp(\theta_t^\top \psi(x, a'))}{\sum_{a' \in \mathcal{A}} \exp(\theta^\top \psi(x, a'))} \right) \\ &\leq (\theta - \theta_t)^\top \psi(x, a) + \log \left(\frac{\sum_{a' \in \mathcal{A}} \exp(\|\theta_t\|_2 \|\psi(x, a')\|_2)}{\sum_{a' \in \mathcal{A}} \exp(\|\theta\|_2 \|\psi(x, a')\|_2)} \right) \\ &\leq (\theta - \theta_t)^\top \psi(x, a) \\ &\leq (\|\theta\|_2 + \|\theta_t\|_2) \|\psi(x, a)\|_2 \\ &\leq 2B \end{aligned}$$

where the first inequality holds from Cauchy-Schwartz, the second inequality holds from upper bounding θ, θ_t, ψ from Assumption 2 and noticing that the second term equates to $\log(1) = 0$. The third inequality from Cauchy-Schwartz and the Triangle inequality, and the last inequality holds again from upper bounding θ, θ_t, ψ . Now, we also have that $r \in \mathcal{F}$. Thus,

$$\begin{aligned} r(x, a) - r(x, a') &= \phi(x, a)^\top \omega - \phi(x, a')^\top \omega \\ &\leq (\|\phi(x, a)\|_2 + \|\phi(x, a')\|_2) \|\omega\|_2 \\ &\leq 2F \end{aligned}$$

where the first inequality holds from Cauchy-Schwartz and Triangle inequality. Now recall the REBEL update 1. Using these facts we have that $|g(z; \theta)| \leq 4B/\eta + 2F$ so $K_g = \sup_{z, \theta} |g(z; \theta)| \leq 4B/\eta + 2F$. Since $\ell(z; \theta) = g(z; \theta)^2$, $K_\ell = K_g^2$. \square

Now we prove that $\ell(z; \theta)$ is $4K_g/\eta$ -Lipschitz in θ .

Lemma 13 (Lipschitz bound on $\ell(z; \theta)$). *$\ell(z; \theta)$ is $\frac{4K_g}{\eta}$ -Lipschitz in θ .*

Proof of Lemma 13. First we compute the gradient $\nabla_\theta g(z; \theta)$. Since we are looking at updates with respect to θ , notice that we have the following:

$$\nabla_\theta g(z; \theta) = \nabla_\theta \left(\frac{1}{\eta} [\log \pi_\theta(a | x) - \log \pi_\theta(a' | x)] \right)$$

Now notice that

$$\begin{aligned} \log \pi_\theta(a | x) - \log \pi_\theta(a' | x) &= \log \left(\exp(\theta^\top \psi(x, a)) \right) - \log \left(\exp(\theta^\top \psi(x, a')) \right) \\ &= \theta^\top (\psi(x, a) - \psi(x, a')) \end{aligned}$$

Thus we find that

$$\nabla_\theta g(z; \theta) = \frac{1}{\eta} (\psi(x, a) - \psi(x, a'))$$

Thus by triangle inequality, $\sup_{x, a} \|\nabla_\theta g(z; \theta)\|_2 \leq 2/\eta$. Now since $\ell(z; \theta) = g(z; \theta)^2$, we have that $\nabla_\theta \ell(z; \theta) = 2g(z; \theta) \nabla_\theta g(z; \theta)$. From Lemma 12, we know that $K_g = \sup_{z, \theta} |g(z; \theta)|$ so we see that $\sup_{x, a} \|\nabla_\theta \ell(z; \theta)\|_2 \leq 4K_g/\eta$. Thus we can conclude that $\ell(z; \theta)$ is $4K_g/\eta$ -Lipschitz in θ . \square

C Proof of "Slow Rate" Wasserstein-DRO-REBEL

First we prove that $h(\theta; \mathbb{P}) = \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$ is strongly convex for any \mathbb{P} .

Lemma 14 (Strong convexity of h). *Let $\ell(z; \theta)$ be the REBEL loss function. Assume that Assumption 3 holds. Then $h(\theta; \mathbb{P}) = \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$ is $2/\eta$ -strongly convex with respect to norm $\|\cdot\|_{\Sigma_{\mathbb{P}}}$ where $\Sigma_{\mathbb{P}} := \mathbb{E}_{(x, a^1, a^2, y) \sim \mathbb{P}} \left[(\psi(x, a^1) - \psi(x, a^2)) (\psi(x, a^1) - \psi(x, a^2))^{\top} \right]$*

Proof of Lemma 14. From Lemma 13, we know that $\nabla_{\theta} \ell(z; \theta) = 2g(z; \theta) \nabla_{\theta} g(z; \theta)$ so we have

$$\nabla_{\theta}^2 \ell(z; \theta) = 2 \nabla_{\theta} g(z; \theta) \nabla_{\theta} g(z; \theta)^{\top}$$

We also know from Lemma 13 that $\nabla_{\theta} g(z; \theta) = \frac{1}{\eta} [\psi(x, a) - \psi(x, a')]$. Thus we find that

$$\nabla_{\theta}^2 \ell(z; \theta) = \frac{2}{\eta^2} (\psi(x, a) - \psi(x, a')) (\psi(x, a) - \psi(x, a'))^{\top}$$

Then taking an expectation under \mathbb{P} , we find that

$$\nabla_{\theta}^2 h(\theta; \mathbb{P}) = \frac{2}{\eta^2} \mathbb{E}_{z \sim \mathbb{P}} \left[(\psi(x, a^1) - \psi(x, a^2)) (\psi(x, a^1) - \psi(x, a^2))^{\top} \right] \succeq \frac{2}{\eta^2} \Sigma_{\mathbb{P}}$$

which follows from the data coverage condition. Now fix $\theta, \theta' \in \Theta$. Let $\Delta = \theta - \theta'$. By the second-order Taylor expansion, there exists $\tilde{\theta}$ on the line segment between θ and θ' such that

$$\ell(z; \theta') - \ell(z; \theta) - \langle \nabla_{\theta} \ell(z; \theta), \Delta \rangle = \frac{1}{2} \Delta^{\top} \nabla_{\tilde{\theta}}^2 \ell(z; \tilde{\theta}) \Delta \geq \frac{\mu}{2} \|\Delta\|_{\Sigma_{\mathbb{P}}}$$

where $\mu = 2/\eta^2$ and $\Sigma_{\mathbb{P}} = (\psi(x, a^1) - \psi(x, a^2)) (\psi(x, a^1) - \psi(x, a^2))^{\top}$. Taking expectations we find that

$$h(\theta'; \mathbb{P}) \geq h(\theta; \mathbb{P}) + \langle \nabla_{\theta} h(z; \theta), \Delta \rangle + \frac{\mu}{2} \|\Delta\|_{\Sigma_{\mathbb{P}}}$$

Thus, h is μ -strongly convex in the $\|\cdot\|_{\Sigma_{\mathbb{P}}}$ norm. □

We now establish strong convexity of $\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^0; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$

Lemma 15 (Strong convexity of $\mathcal{L}^{\mathcal{W}_p}$). *Let $\ell(z; \theta)$ be the REBEL loss function. Then $\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon}(\mathbb{P}^0; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)]$ is $2\lambda/\eta^2$ -strongly convex with respect to Euclidean norm $\|\cdot\|_2$ where λ is the regularity parameter from Assumption 3.*

Proof of Lemma 19. In Lemma 14, we proved strong convexity of h . By Lemma 5, for $\theta, \theta' \in \Theta$ and $\alpha \in [0, 1]$, this is equivalent to

$$h(\alpha\theta + (1 - \alpha)\theta'; \mathbb{P}) \leq \alpha h(\theta; \mathbb{P}) + (1 - \alpha) h(\theta'; \mathbb{P}) - \frac{\mu}{2} \alpha(1 - \alpha) \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2$$

Taking the supremum over \mathbb{P} preserves the convex combination and the negative quadratic term so we get

$$\begin{aligned}
\mathcal{L}^{\mathcal{W}_p}(\alpha\theta + (1-\alpha)\theta'; \varepsilon) &= \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} h(\alpha\theta + (1-\alpha)\theta'; \mathbb{P}) \\
&\leq \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \left[\alpha h(\theta; \mathbb{P}) + (1-\alpha)h(\theta'; \mathbb{P}) - \frac{\mu}{2}\alpha(1-\alpha)\|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2 \right] \\
&\leq \alpha \mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\mathcal{W}_p}(\theta'; \varepsilon) - \frac{\mu}{2}\alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2 \\
&\leq \alpha \mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\mathcal{W}_p}(\theta'; \varepsilon) - \frac{\mu}{2}\alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \lambda_{\min}(\Sigma_{\mathbb{P}}) \|\theta - \theta'\|_2^2 \\
&\leq \alpha \mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\mathcal{W}_p}(\theta'; \varepsilon) - \frac{\mu\lambda}{2}\alpha(1-\alpha)\|\theta - \theta'\|_2^2
\end{aligned}$$

where the second inequality holds from $\sup_x (f(x) + g(x)) \leq \sup_x f(x) + \sup_x g(x)$, the third inequality holds by the fact that $\Sigma_{\mathbb{P}} \succeq \lambda_{\min}(\Sigma_{\mathbb{P}}) I$, and the last inequality holds from Assumption 3. Thus we conclude that $\mathcal{L}^{\mathcal{W}_p}$ is $\mu\lambda$ -strongly convex in the $\|\cdot\|_2$ norm. \square

We are now ready to prove the "slow rate" estimation error of Wasserstein-DRO-REBEL.

Proof of Theorem 1. By strong duality for Wasserstein DRO [2], for fixed θ we have

$$\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] = \inf_{\Delta \geq 0} \{ \delta \varepsilon^p - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)] \}$$

where $\ell_\Delta(z; \theta) = \inf_{z' \in \mathcal{Z}} \{ \Delta d^p(z, z') - \ell(z'; \theta) \}$ where d is the metric used to define the type- p Wasserstein distance. Then notice that

$$\begin{aligned}
|\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon)| &= \left| \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}_n^\circ; \mathcal{W}_p)} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] \right| \\
&= \left| \inf_{\Delta \geq 0} \{ \Delta \varepsilon^p - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)] \} - \inf_{\Delta \geq 0} \{ \Delta \varepsilon^p - \mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell_\Delta(z; \theta)] \} \right| \\
&\leq \sup_{\Delta \geq 0} |\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell_\Delta(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)]|
\end{aligned}$$

where the first equality holds from strong duality and the last inequality holds from $\inf_x f(x) - \inf_x g(x) \leq \sup_x |f(x) - g(x)|$. From Lemma 12, we showed that $\ell(z; \theta) \in [0, K_\ell]$. Now notice that

$$\begin{aligned}
l_\Delta(z; \theta) &= \inf_{z' \in \mathcal{Z}} \{ \Delta d^p(z, z') - \ell(z'; \theta) \} \leq \inf_{z' \in \mathcal{Z}} \{ \Delta d^p(z, z') \} = 0 \\
\ell_\Delta(z; \theta) &= \inf_{z' \in \mathcal{Z}} \{ \Delta d^p(z, z') - \ell(z'; \theta) \} \geq \inf_{z' \in \mathcal{Z}} \{ \Delta d^p(z, z') - K_\ell \} \geq -K_\ell
\end{aligned}$$

Thus, $\ell_\Delta \in [-K_\ell, 0]$. Since ℓ_Δ is bounded and $z \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_n^\circ$, we have by Hoeffding's inequality (by Lemma 11)

$$\mathbb{P}(|\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell_\Delta(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell_\Delta(z; \theta)]| \geq \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2}{K_\ell^2}\right)$$

\square

Since K_ℓ does not depend on Δ , this bound is uniform and thus does not require any uniform covering argument. By picking δ to be the right hand side, we find that with probability at least $1 - \delta$

$$|\mathcal{L}^{\mathcal{W}_p}(\theta; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon)| \leq K_\ell \sqrt{\frac{\log(2/\delta)}{2n}}$$

Now we have that

$$\begin{aligned}
& \mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon) \\
&= \mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) + \mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon) + \mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon) \\
&\leq |\mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon)| + |\mathcal{L}_n^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}_n^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon)| \\
&\leq K_\ell \sqrt{\frac{2 \log(2/\delta)}{n}}
\end{aligned}$$

where the first inequality holds from the fact that $\hat{\theta}_n^{\mathcal{W}_p} \in \arg \min_{\theta \in \Theta} \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon)$. Now from Lemma 6 and Lemma 19, we have that

$$\frac{\lambda}{\eta^2} \|\theta^{\mathcal{W}_p} - \hat{\theta}_n^{\mathcal{W}_p}\|^2 \leq \mathcal{L}^{\mathcal{W}_p}(\theta^{\mathcal{W}_p}; \varepsilon) - \mathcal{L}^{\mathcal{W}_p}(\hat{\theta}_n^{\mathcal{W}_p}; \varepsilon)$$

Thus with probability at least $1 - \delta$, we conclude that

$$\|\theta^{\mathcal{W}_p} - \hat{\theta}_n^{\mathcal{W}_p}\|^2 \leq \frac{\eta^2 K_g^2}{\lambda} \sqrt{\frac{2 \log(2/\delta)}{n}}$$

D Proof of "Slow Rate" KL-DRO-REBEL

Before we prove the necessary results to get the "slow rate" for KL-DRO-REBEL, we need to make an assumption on the loss functions $\ell(\cdot; \theta)$, $\theta \in \Theta$. Note that this assumption is only used in proving the dual reformulation of the KL-DRO-REBEL objective.

Assumption 4. We assume that $\ell(z; \theta) \leq L$ for all $\theta \in \Theta$. That is, the loss function is upper bounded by L . In addition, we also assume that Θ permits a uniform upper bound on λ_θ . That is, we assume that

$$\sup_{\theta \in \Theta} \lambda_\theta < \bar{\lambda}.$$

We state the following dual reformulation result. The proof of this reformulation can be found in [Xu et al., 2025], Appendix C:

Lemma 16 (Dual reformulation of KL-DRO-REBEL). *Let $\ell(z; \theta)$ be the REBEL loss. The KL-DRO-REBEL loss function admits the following dual reformulation:*

$$\mathcal{L}^{\text{KL}}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)] = \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left\{ \lambda \varepsilon + \lambda \log \left(\mathbb{E}_{z \sim \mathbb{P}^\circ} \left[\exp \left(\frac{\ell(z; \theta)}{\lambda} \right) \right] \right) \right\},$$

where $0 < \underline{\lambda} < \bar{\lambda} < \infty$ are constants.

We will now establish strong convexity of $\mathcal{L}^{\text{KL}}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)]$. This proof will essentially be the same as Lemma 19

Lemma 17 (Strong convexity of \mathcal{L}^{KL}). *Let $\ell(z; \theta)$ be the REBEL loss function. Then $\mathcal{L}^{\text{KL}}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)]$ is $2\lambda/\eta^2$ -strongly convex with respect to Euclidean norm $\|\cdot\|_2$ where λ is the regularity parameter from Assumption 3.*

Proof of Lemma 19. In Lemma 14, we proved strong convexity of h . By Lemma 5, for $\theta, \theta' \in \Theta$ and $\alpha \in [0, 1]$, this is equivalent to

$$h(\alpha\theta + (1 - \alpha)\theta'; \mathbb{P}) \leq \alpha h(\theta; \mathbb{P}) + (1 - \alpha)h(\theta'; \mathbb{P}) - \frac{\mu}{2} \alpha(1 - \alpha) \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2$$

Taking the supremum over \mathbb{P} preserves the convex combination and the negative quadratic term so we get

$$\begin{aligned}
\mathcal{L}^{\text{KL}}(\alpha\theta + (1-\alpha)\theta'; \varepsilon) &= \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} h(\alpha\theta + (1-\alpha)\theta'; \mathbb{P}) \\
&\leq \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \left[\alpha h(\theta; \mathbb{P}) + (1-\alpha) h(\theta'; \mathbb{P}) - \frac{\mu}{2} \alpha(1-\alpha) \|\theta - \theta'\|_{\Sigma_\mathbb{P}}^2 \right] \\
&\leq \alpha \mathcal{L}^{\text{KL}}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\text{KL}}(\theta'; \varepsilon) - \frac{\mu}{2} \alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \|\theta - \theta'\|_{\Sigma_\mathbb{P}}^2 \\
&\leq \alpha \mathcal{L}^{\text{KL}}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\text{KL}}(\theta'; \varepsilon) - \frac{\mu}{2} \alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \lambda_{\min}(\Sigma_\mathbb{P}) \|\theta - \theta'\|_2^2 \\
&\leq \alpha \mathcal{L}^{\text{KL}}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\text{KL}}(\theta'; \varepsilon) - \frac{\mu\lambda}{2} \alpha(1-\alpha) \|\theta - \theta'\|_2^2
\end{aligned}$$

where the second inequality holds from $\sup_x (f(x) + g(x)) \leq \sup_x f(x) + \sup_x g(x)$, the third inequality holds by the fact that $\Sigma_\mathbb{P} \succeq \lambda_{\min}(\Sigma_\mathbb{P}) I$, and the last inequality holds from Assumption 3. Thus we conclude that \mathcal{L}^{KL} is $\mu\lambda$ -strongly convex in the $\|\cdot\|_2$ norm. \square

We are now ready to prove the "slow rate" estimation error of KL-DRO-REBEL.

Proof of Theorem 2. By the strong duality result for KL-DRO [16], we have for fixed θ

$$\mathcal{L}^{\text{KL}}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] = \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \{ \lambda \varepsilon + \lambda \log(\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]) \},$$

where $j(z, \lambda; \theta) = \exp\left(\frac{\ell(z; \theta)}{\lambda}\right)$. Then we have

$$\begin{aligned}
|\mathcal{L}^{\text{KL}}(\theta; \varepsilon) - \mathcal{L}_n^{\text{KL}}(\theta; \varepsilon)| &= \left| \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}_n^\circ; \text{KL})} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] \right| \\
&= \left| \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \{ \lambda \varepsilon + \lambda \log(\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]) \} - \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \{ \lambda \varepsilon + \lambda \log(\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)]) \} \right| \\
&\leq \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left| \lambda \log(\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)]) - \lambda \log(\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]) \right| \\
&= \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \lambda \left| \log \left(\frac{\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)]}{\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]} \right) \right| \\
&= \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \lambda \left| \log \left(\frac{\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]}{\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]} + 1 \right) \right| \\
&\leq \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \lambda \left| \frac{\mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]}{\mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)]} \right| \\
&\leq \bar{\lambda} \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left| \mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)] \right|
\end{aligned}$$

where the first equality holds from strong duality, the first inequality holds from $\inf_x f(x) - \inf_x g(x) \leq \sup_x |f(x) - g(x)|$, the second inequality holds from $|\log(1+x)| \leq |x| \forall x \geq 0$, and the last inequality holds from $\ell(z; \theta) \geq 0$. By Hoeffding's inequality (by Lemma 11)

$$\mathbb{P} \left(\left| \mathbb{E}_{z \sim \mathbb{P}_n^\circ} [j(z, \lambda; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [j(z, \lambda; \theta)] \right| \geq \epsilon \right) \leq 2 \exp \left(- \frac{2n\epsilon^2}{\exp(K_g^2/\bar{\lambda})} \right)$$

Again since K_g and $\bar{\lambda}$ are independent of λ , we do not need a uniform covering argument over λ . Picking δ to be the right side, we find that with probability at least $1 - \delta$

$$|\mathcal{L}^{\text{KL}}(\theta; \varepsilon) - \mathcal{L}_n^{\text{KL}}(\theta; \varepsilon)| \leq \bar{\lambda} \sqrt{\frac{\exp(K_g^2/\bar{\lambda}) \log(2/\delta)}{2n}}$$

Now we have that

$$\begin{aligned}
& \mathcal{L}^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) - \mathcal{L}^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) \\
&= \mathcal{L}^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) - \mathcal{L}_n^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) + \mathcal{L}_n^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) - \mathcal{L}_n^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) + \mathcal{L}_n^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) - \mathcal{L}^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) \\
&\leq |\mathcal{L}^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) - \mathcal{L}_n^{\text{KL}}(\theta^{\text{KL}}; \varepsilon)| + |\mathcal{L}_n^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) - \mathcal{L}^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon)| \\
&\leq \bar{\lambda} \sqrt{\frac{2 \exp(K_g^2/\bar{\lambda}) \log(2/\delta)}{n}}
\end{aligned}$$

where the first inequality holds from the fact that $\hat{\theta}_n^{\mathcal{W}_p} \in \arg \min_{\theta \in \Theta} \mathcal{L}_n^{\mathcal{W}_p}(\theta; \varepsilon)$. Now from Lemma 6 and Lemma 16, we have that

$$\frac{\lambda}{\eta^2} \|\theta^{\text{KL}} - \hat{\theta}_n^{\text{KL}}\|^2 \leq \left| \mathcal{L}^{\text{KL}}(\theta^{\text{KL}}; \varepsilon) - \mathcal{L}^{\text{KL}}(\hat{\theta}_n^{\text{KL}}; \varepsilon) \right|$$

Thus with probability at least $1 - \delta$, we conclude that

$$\|\theta^{\text{KL}} - \hat{\theta}_n^{\text{KL}}\|^2 \leq \frac{\eta^2}{\lambda} \sqrt{\frac{2\bar{\lambda}^2 \exp(K_g^2/\bar{\lambda}) \log(2/\delta)}{n}}$$

□

E Proof of "Slow Rate" χ^2 -DRO-REBEL

We first state the following dual reformulation for χ^2 -DRO

Lemma 18 (Dual reformulation of χ^2 -DRO-REBEL). *Let $\ell(z; \theta)$ be the REBEL loss. The χ^2 -DRO-REBEL objective admits the dual form*

$$\begin{aligned}
\mathcal{L}^{\chi^2}(\theta; \varepsilon) &= \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)] \\
&= \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left\{ \lambda \varepsilon + \mathbb{E}_{z \sim \mathbb{P}^\circ}[\ell(z; \theta)] - 2\lambda + \frac{1}{4\lambda} \mathbb{E}_{z \sim \mathbb{P}^\circ}[(\ell(z; \theta) - \mathbb{E}_{\mathbb{P}^\circ}[\ell(z; \theta)] + 2\lambda)^2] \right\}
\end{aligned}$$

where $0 < \underline{\lambda} < \bar{\lambda} < \infty$ are chosen so that the infimum is attained. Equivalently, defining $\mu = \mathbb{E}_{\mathbb{P}^\circ}[\ell(z; \theta)]$ and $\sigma^2 = \text{Var}_{\mathbb{P}^\circ}(\ell(z; \theta))$,

$$\mathcal{L}^{\chi^2}(\theta; \varepsilon) = \mu + \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left\{ (\varepsilon - 1)\lambda + \frac{\sigma^2}{4\lambda} \right\}.$$

We again will prove strong convexity for χ^2

Lemma 19 (Strong convexity of \mathcal{L}^{χ^2}). *Let $l(z; \theta)$ be the REBEL loss function. Then $\mathcal{L}^{\chi^2}(\theta; \varepsilon) = \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} \mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)]$ is $2\lambda/\eta$ -strongly convex with respect to Euclidean norm $\|\cdot\|_2$ where λ is the regularity parameter from Assumption 3.*

Proof of Lemma 19. In Lemma 14, we proved strong convexity of h . By Lemma 5, for $\theta, \theta' \in \Theta$ and $\alpha \in [0, 1]$, this is equivalent to

$$h(\alpha\theta + (1 - \alpha)\theta'; \mathbb{P}) \leq \alpha h(\theta; \mathbb{P}) + (1 - \alpha)h(\theta'; \mathbb{P}) - \frac{\mu}{2} \alpha(1 - \alpha) \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2$$

Taking the supremum over \mathbb{P} preserves the convex combination and the negative quadratic term so we get

$$\begin{aligned}
\mathcal{L}^{\chi^2}(\alpha\theta + (1-\alpha)\theta'; \varepsilon) &= \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} h(\alpha\theta + (1-\alpha)\theta'; \mathbb{P}) \\
&\leq \sup_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} \left[\alpha h(\theta; \mathbb{P}) + (1-\alpha) h(\theta'; \mathbb{P}) - \frac{\mu}{2} \alpha(1-\alpha) \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2 \right] \\
&\leq \alpha \mathcal{L}^{\chi^2}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\chi^2}(\theta'; \varepsilon) - \frac{\mu}{2} \alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} \|\theta - \theta'\|_{\Sigma_{\mathbb{P}}}^2 \\
&\leq \alpha \mathcal{L}^{\chi^2}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\chi^2}(\theta'; \varepsilon) - \frac{\mu}{2} \alpha(1-\alpha) \inf_{\mathbb{P} \in \mathcal{B}_\varepsilon(\mathbb{P}^\circ; \chi^2)} \lambda_{\min}(\Sigma_{\mathbb{P}}) \|\theta - \theta'\|_2^2 \\
&\leq \alpha \mathcal{L}^{\chi^2}(\theta; \varepsilon) + (1-\alpha) \mathcal{L}^{\chi^2}(\theta'; \varepsilon) - \frac{\mu\lambda}{2} \alpha(1-\alpha) \|\theta - \theta'\|_2^2
\end{aligned}$$

where the second inequality holds from $\sup_x (f(x) + g(x)) \leq \sup_x f(x) + \sup_x g(x)$, the third inequality holds by the fact that $\Sigma_{\mathbb{P}} \succeq \lambda_{\min}(\Sigma_{\mathbb{P}}) I$, and the last inequality holds from Assumption 3. Thus we conclude that \mathcal{L}^{χ^2} is $\mu\lambda$ -strongly convex in the $\|\cdot\|_2$ norm. \square

We now prove the "slow rate" estimation error of χ^2 -DRO-REBEL

Proof of Theorem 3. Let $\theta^{\chi^2} \in \arg \min_{\theta} \mathcal{L}^{\chi^2}(\theta; \varepsilon)$ and $\hat{\theta}_n^{\chi^2} \in \arg \min_{\theta} \mathcal{L}_n^{\chi^2}(\theta; \varepsilon)$. By the dual reformulation (Lemma 18), for any fixed θ

$$\mathcal{L}^{\chi^2}(\theta; \varepsilon) = \mu + \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left\{ (\varepsilon - 1) \lambda + \frac{\sigma^2}{4\lambda} \right\},$$

and similarly

$$\mathcal{L}_n^{\chi^2}(\theta; \varepsilon) = \mu_n + \inf_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} \left\{ (\varepsilon - 1) \lambda + \frac{\sigma_n^2}{4\lambda} \right\},$$

where $\mu = \mathbb{E}_{\mathbb{P}^\circ}[\ell(z; \theta)]$, $\mu_n = \mathbb{E}_{\mathbb{P}_n^\circ}[\ell(z; \theta)]$, $\sigma^2 = \text{Var}_{\mathbb{P}^\circ}(\ell(z; \theta))$, and $\sigma_n^2 = \text{Var}_{\mathbb{P}_n^\circ}(\ell(z; \theta))$. Using the dual reformulation we have that

$$\begin{aligned}
|\mathcal{L}^{\chi^2}(\theta; \varepsilon) - \mathcal{L}_n^{\chi^2}(\theta; \varepsilon)| &= |\mu - \mu_n| + \left| \inf_{\lambda} g(\lambda) - \inf_{\lambda} g_n(\lambda) \right| \\
&\leq |\mu - \mu_n| + \sup_{\lambda \in [\underline{\lambda}, \bar{\lambda}]} |g(\lambda) - g_n(\lambda)|,
\end{aligned}$$

where $g(\lambda) = (\varepsilon - 1)\lambda + \frac{\sigma^2}{4\lambda}$ and $g_n(\lambda) = (\varepsilon - 1)\lambda + \frac{\sigma_n^2}{4\lambda}$. The second inequality in the argument above follows by the triangle inequality. From Lemma 12, we have that $\ell(z; \theta) \in [0, K_\ell]$. Then by Hoeffding's inequality (Lemma 11), we have that

$$\Pr(|\mu - \mu_n| \geq \delta) \leq 2 \exp\left(-\frac{2n\delta^2}{K_\ell^2}\right),$$

and

$$\Pr(|\sigma^2 - \sigma_n^2| \geq \delta') \leq 2 \exp\left(-\frac{2n\delta'^2}{K_\ell^4}\right).$$

Thus with probability at least $1 - 2 \exp(-2n\alpha^2/K_\ell^4)$, choosing $\delta, \delta' = K_\ell^2 \sqrt{\frac{\log(2/\delta)}{2n}}$, we get

$$|\mu - \mu_n| \leq K_\ell \sqrt{\frac{\log(2/\delta)}{2n}},$$

and

$$|\sigma^2 - \sigma_n^2| \leq K_\ell^2 \sqrt{\frac{\log(2/\delta)}{2n}}.$$

Therefore

$$\sup_{\lambda} |g(\lambda) - g_n(\lambda)| \leq \frac{|\sigma^2 - \sigma_n^2|}{4\bar{\lambda}} \leq \frac{K_\ell^2}{4\bar{\lambda}} \sqrt{\frac{\log(2/\delta)}{2n}},$$

and overall

$$|\mathcal{L}^{\chi^2}(\theta; \varepsilon) - \mathcal{L}_n^{\chi^2}(\theta; \varepsilon)| \leq K_\ell \left(1 + K_\ell/4\bar{\lambda}\right) \sqrt{\frac{\log(2/\delta)}{2n}}$$

Now set $\theta = \theta^{\chi^2}$ and $\theta = \hat{\theta}_n^{\chi^2}$ in turn, and use the “three-term” decomposition (as in the Wasserstein and KL proof) to conclude

$$\mathcal{L}^{\chi^2}(\theta^{\chi^2}; \varepsilon) - \mathcal{L}^{\chi^2}(\hat{\theta}_n^{\chi^2}; \varepsilon) \leq 2|\mathcal{L}^{\chi^2}(\theta) - \mathcal{L}_n^{\chi^2}(\theta)| \leq K_\ell \left(1 + K_\ell/4\bar{\lambda}\right) \sqrt{\frac{2\log(2/\delta)}{n}}$$

Finally, by strong convexity of \mathcal{L}^{χ^2} (cf. Lemma 6),

$$\frac{\lambda}{\eta^2} \|\theta^{\chi^2} - \hat{\theta}_n^{\chi^2}\|^2 \leq \mathcal{L}^{\chi^2}(\theta^{\chi^2}; \varepsilon) - \mathcal{L}^{\chi^2}(\hat{\theta}_n^{\chi^2}; \varepsilon),$$

so with probability at least $1 - \delta$,

$$\|\theta^{\chi^2} - \hat{\theta}_n^{\chi^2}\|^2 \leq \frac{\eta^2 K_g^2}{\lambda} \left(1 + \frac{K_g^2}{4\bar{\lambda}}\right) \sqrt{\frac{2\log(2/\delta)}{n}}$$

as claimed. \square

F Proof of "Master Theorem" for Parametric $n^{-1/2}$ rates

Proof of Theorem 4. Let $M(\theta) = \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell(z; \theta)]$ and $M_n(\theta) = \mathbb{E}_{z \sim \mathbb{P}_n^\circ} [\ell(z; \theta)]$. First notice that we can do a loss decomposition as follows: for fixed $\theta \in \Theta$

$$\begin{aligned} |\mathcal{L}^D(\theta; \varepsilon_n) - \mathcal{L}_n^D(\theta; \varepsilon_n)| &\leq |\mathcal{L}^D(\theta; \varepsilon_n) - M(\theta)| + |M(\theta) - M_n(\theta)| + |M_n(\theta) - \mathcal{L}_n^D(\theta; \varepsilon_n)| \\ &\leq 2\Delta_n + \sup_{\|\theta - \theta^*\|_2 \leq r} |M(\theta) - M_n(\theta)| \end{aligned}$$

where the first inequality holds from triangle inequality and the second holds from the assumption that $|\mathcal{L}^D(\theta) - \mathbb{E}_{\mathbb{P}}[\ell(z; \theta)]| \leq \Delta_n$ and $|\mathcal{L}_n^D(\theta; \varepsilon_n) - \mathbb{E}_{\mathbb{P}_n^\circ}[\ell(z; \theta)]| \leq \Delta_n$. Now let us define the following function class

$$\mathcal{F}_r = \{f_\theta(z) = \ell(z; \theta) - \ell(z; \theta^*) \mid \|\theta - \theta^*\|_2 \leq r\}$$

Then bounding $\sup_{\|\theta - \theta^*\|_2 \leq r} |M(\theta) - M_n(\theta)|$ is equivalent to bounding $\sup_{f \in \mathcal{F}_r} |P_n f - P f|$ where $P_n f - P f$ corresponds to the empirical process $M(\theta) - M_n(\theta)$. Now by Lemma 8 (Symmetrization) taking $\Phi = I$ and using the notation $\mathcal{R}_n(\mathcal{F}_r) = \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}_r} \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right]$ called the Rademacher complexity, we have

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}_r} |P_n f - P f| \right] \leq 2\mathbb{E} [\mathcal{R}_n(\mathcal{F}_r)]$$

where the expectation is taken with respect to the data z in the expression above whereas in the Rademacher complexity, we condition on the data z and take the expectation over $\sigma_i \stackrel{\text{i.i.d.}}{\sim} \text{Unif}\{-1, 1\}$. Assume that $\ell(z; \theta)$ is L_g -Lipschitz in θ . A standard covering argument gives us

$$N(\epsilon, \mathcal{F}_r, L_2(P)) \leq \left(\frac{3L_g r}{\epsilon} \right)^d$$

Corollary 5 (Dudley’s entropy integral) then gives us

$$\mathbb{E} [\mathcal{R}_n(\mathcal{F}_r)] \leq \frac{12}{\sqrt{n}} \int_0^{L_g r} \sqrt{\log N(\epsilon, \mathcal{F}_r, L_2(P))} d\epsilon \leq c_0 L_g r \sqrt{\frac{d}{n}}$$

where $c_0 = \int_0^1 \sqrt{\log(3/u)} du < \infty$ is an absolute constant. Thus we have that

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}_r} |P_n f - P f| \right] \leq 2c_0 L_g r \sqrt{\frac{d}{n}}$$

To upgrade the above expectation bound into a high-probability bound, we will apply Bousquet’s inequality (Theorem 7). First notice that since $\ell(z; \theta)$ is L_g -Lipschitz in θ , we have that

$$|f_\theta(z)| = |\ell(z; \theta) - \ell(z; \theta^*)| \leq L_g \|\theta - \theta^*\|_2 \leq L_g r$$

Thus $f_\theta^2(z) \leq L_g^2 r^2$. Thus we have that $\text{Var}_{z \sim \mathbb{P}}(f_\theta(z)) \leq L_g^2 r^2$. Since each $f_\theta \in \mathcal{F}_r$ is uniformly bounded by $2K_\ell$ by assumption and has bounded variance, then by applying Theorem 7 to Rademacher averages of a class, we find that there exists constants $c_1, c_2 > 0$ such that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$

$$\sup_{f \in \mathcal{F}_r} |P_n f - P f| \leq 2\mathbb{E}[\mathcal{R}_n(\mathcal{F}_r)] + c_1 L_g r \sqrt{\frac{\log(1/\delta)}{n}} + c_2 \frac{K_\ell}{n} \log(1/\delta)$$

Plugging in the Dudley integral bound, we get

$$\sup_{f \in \mathcal{F}_r} |P_n f - P f| \leq 2c_0 L_g r \sqrt{\frac{d}{n}} + c_1 L_g r \sqrt{\frac{\log(1/\delta)}{n}} + c_2 \frac{K_\ell}{n} \log(1/\delta)$$

We use a **localization argument** to avoid assuming $\hat{\theta}_n$ is in a fixed ball from the start. Let $R_{\max} = 2B$ be the maximum possible distance in the parameter space (since $\theta \in \Theta$ and $\|\theta\|_2 \leq B$). Define a sequence of radii $r_k = 2^k r_0$ for $k = 0, 1, \dots, K_{\max} = \lceil \log_2(R_{\max}/r_0) \rceil$, where r_0 is a small positive constant and $2^{K_{\max}} r_0 \geq R_{\max}$. This creates a series of nested balls $B(\theta^*, r_1) \subset \dots \subset B(\theta^*, r_{K_{\max}})$ covering Θ . For each k , define an event \mathcal{E}_k that the uniform convergence bound holds for the ball $B(\theta^*, r_k)$:

$$\mathcal{E}_k = \left\{ \sup_{\|\theta - \theta^*\|_2 \leq r_k} |\mathcal{L}_n^D(\theta; \varepsilon_n) - \mathcal{L}^D(\theta; \varepsilon_n)| \leq \psi_n(r_k, \delta_k) \right\}$$

where

$$\psi_n(r, \delta) = 2\Delta_n + 2c_0 L_g r \sqrt{\frac{d}{n}} + c_1 L_g r \sqrt{\frac{\log(1/\delta)}{n}} + c_2 \frac{K_\ell}{n} \log(1/\delta)$$

We set $\delta_k = \delta/(2^{k+1})$, so $\sum_{k=0}^{K_{\max}} \delta_k < \delta$. By a union bound, the event $\mathcal{E} = \bigcap_{k=0}^{K_{\max}} \mathcal{E}_k$ holds with probability at least $1 - \sum \delta_k > 1 - \delta$. On this event \mathcal{E} , the uniform convergence bound holds for *any* θ within *any* $B(\theta^*, r_k)$. Let $\tilde{r} = \|\hat{\theta}_n - \theta^*\|_2$. Since $\hat{\theta}_n \in \Theta$, we know $\tilde{r} \leq R_{\max}$. Let k^* be the smallest integer such that $\tilde{r} \leq r_{k^*}$. By this definition, $\hat{\theta}_n \in B(\theta^*, r_{k^*})$. If $k^* > 0$, then $r_{k^*-1} < \tilde{r} \leq r_{k^*}$. Since $r_{k^*} = 2 \cdot r_{k^*-1}$ (for a dyadic sequence), this implies $r_{k^*} < 2\tilde{r}$. If $k^* = 0$, then $\tilde{r} \leq r_0$, and we use r_0 as the effective radius in the bound. In either case, we can use an upper bound of $2\tilde{r}$ for r_{k^*} (or r_0 if \tilde{r} is extremely small) in the linear term of ψ_n for large n . Since $\hat{\theta}_n$ is the empirical minimizer:

$$\mathcal{L}_n^D(\hat{\theta}_n; \varepsilon_n) \leq \mathcal{L}_n^D(\theta^*; \varepsilon_n)$$

Now we use the uniform convergence bound on the event \mathcal{E} . Since $\hat{\theta}_n \in B(\theta^*, r_{k^*})$, and $\theta^* \in B(\theta^*, r_{k^*})$, we can apply the uniform convergence bound with radius r_{k^*} (and total error probability δ):

$$\begin{aligned} \mathcal{L}_n^D(\hat{\theta}_n; \varepsilon_n) &\geq \mathcal{L}^D(\hat{\theta}_n; \varepsilon_n) - \psi_n(r_{k^*}, \delta) \\ \mathcal{L}_n^D(\theta^*; \varepsilon_n) &\leq \mathcal{L}^D(\theta^*; \varepsilon_n) + \psi_n(r_{k^*}, \delta) \end{aligned}$$

Combining these with the strong convexity of $\mathcal{L}^D(\theta; \varepsilon_n)$ around θ^* :

$$\begin{aligned} \mathcal{L}^D(\theta^*; \varepsilon_n) + \frac{\alpha}{2} \|\hat{\theta}_n - \theta^*\|_2^2 - \psi_n(r_{k^*}, \delta') &\leq \mathcal{L}_n^D(\hat{\theta}_n; \varepsilon_n) \\ &\leq \mathcal{L}_n^D(\theta^*; \varepsilon_n) \\ &\leq \mathcal{L}^D(\theta^*; \varepsilon_n) + \psi_n(r_{k^*}, \delta) \end{aligned}$$

Subtracting $\mathcal{L}^D(\theta^*; \varepsilon_n)$ from all parts and simplifying, we get:

$$\frac{\alpha}{2} \|\hat{\theta}_n - \theta^*\|_2^2 \leq 2\psi_n(r_{k^*}, \delta)$$

Substituting $r_{k^*} \leq 2\tilde{r}$ (where $\tilde{r} = \|\hat{\theta}_n - \theta^*\|_2$) into the definition of $\psi_n(r_{k^*}, \delta')$ for the linear term in r :

$$\frac{\alpha}{2} \tilde{r}^2 \leq 2 \left(2\Delta_n + \left(2c_0 L_g \sqrt{\frac{d}{n}} + c_1 L_g \sqrt{\frac{\log(1/\delta)}{n}} \right) (2\tilde{r}) + c_2 \frac{K_\ell}{n} \log(1/\delta) \right)$$

Rearranging into a quadratic inequality of the form $A\tilde{r}^2 - B\tilde{r} - C \leq 0$:

$$\frac{\alpha}{2} \tilde{r}^2 - \left(4c_0 L_g \sqrt{\frac{d}{n}} + 2c_1 L_g \sqrt{\frac{\log(1/\delta)}{n}} \right) \tilde{r} - \left(4\Delta_n + 2c_2 \frac{K_\ell}{n} \log(1/\delta) \right) \leq 0$$

For such a quadratic $Ax^2 - Bx - C \leq 0$ with $A, B, C > 0$, the solutions are bounded by $x \leq \frac{B + \sqrt{B^2 + 4AC}}{2A}$. Applying this to \tilde{r} :

$$\begin{aligned}\tilde{r} &\leq \left(\frac{4c_0 L_g}{\alpha} \sqrt{\frac{d}{n}} + \frac{2c_1 L_g}{\alpha} \sqrt{\frac{\log(1/\delta)}{n}} \right) + \sqrt{\frac{1}{\alpha^2} \left(4c_0 L_g \sqrt{\frac{d}{n}} + 2c_1 L_g \sqrt{\frac{\log(1/\delta)}{n}} \right)^2 + \frac{4}{\alpha} \left(4\Delta_n + 2c_2 \frac{K_\ell}{n} \log(1/\delta) \right)} \\ &\leq \left(\frac{4c_0 L_g}{\alpha} \sqrt{\frac{d}{n}} + \frac{2c_1 L_g}{\alpha} \sqrt{\frac{\log(1/\delta)}{n}} \right) + \frac{1}{\alpha} \left(4c_0 L_g \sqrt{\frac{d}{n}} + 2c_1 L_g \sqrt{\frac{\log(1/\delta)}{n}} \right) + \sqrt{\frac{8\Delta_n}{\alpha}} + \sqrt{\frac{4c_2 K_\ell}{\alpha n} \log(1/\delta)}\end{aligned}$$

where the second inequality holds from $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$. Combining terms, with probability at least $1 - \delta$, the estimation error is bounded by:

$$\|\hat{\theta}_n - \theta^*\|_2 \leq \frac{16c_0 L_g}{\alpha} \sqrt{\frac{d}{n}} + \frac{8c_1 L_g}{\alpha} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8\Delta_n}{\alpha}} + \sqrt{\frac{4c_2 K_\ell \log(1/\delta)}{\alpha n}}$$

Assuming the DRO approximation error Δ_n (which reflects the inherent gap between the true and nominal risks in the ambiguity set) decays as $O(n^{-1})$, every term on the right-hand side is of order $O(n^{-1/2})$, which concludes the proof. \square

For all proofs that follow, note that by the Master Theorem (Theorem 4), it suffices to show for the metric we are dealing with that we have a dual "remainder" term as follows:

$$|\mathcal{L}^D(\theta; \varepsilon_n) - \mathbb{E}_{\mathbb{P}^\circ}[\ell(z; \theta)]| \leq \Delta_n.$$

Furthermore, the exact same argument applies to the empirical counterpart, i.e.,

$$|\mathcal{L}_n^D(\theta; \varepsilon_n) - \mathbb{E}_{\mathbb{P}_n^\circ}[\ell(z; \theta)]| \leq \Delta_n,$$

by simply replacing \mathbb{P}° with \mathbb{P}_n° in the derivations. This holds because the properties of the loss function and the definitions of the divergences are universally applicable to any probability measure. In all cases, we will show that $\Delta_n = O(n^{-1})$ by choosing ε_n appropriately.

G Proof of "Fast Rate" Wasserstein-DRO-REBEL

Before we prove the necessary results to get the "slow rate" for KL-DRO-REBEL, we need to make an assumption on the loss functions $\ell(\cdot; \theta)$, $\theta \in \Theta$. Note that this assumption is only used in proving the dual "remainder" term holds:

Assumption 5. For the Wasserstein-DRO-REBEL objective, we assume that the pointwise loss function $\ell(z; \theta)$ is $\mathbf{L}_{\ell, z}$ -Lipschitz with respect to its data argument $z = (x, a^1, a^2)$ for all $\theta \in \Theta$. That is, there exists a constant $L_{\ell, z} \geq 0$ such that for any $z_1, z_2 \in \mathcal{Z}$,

$$|\ell(z_1; \theta) - \ell(z_2; \theta)| \leq L_{\ell, z} d(z_1, z_2)$$

where $d(\cdot, \cdot)$ is the metric corresponding to the type- p Wasserstein distance used to define the ambiguity set.

Proof of Corollary 1. First, recall the type- p Wasserstein distance. The type- p ($p \in [1, \infty)$) Wasserstein distance between two distributions $\mathbb{P}, \mathbb{Q} \in \mathcal{M}(\Xi)$ is defined as

$$\mathcal{W}_p(\mathbb{P}, \mathbb{Q}) = \left(\inf_{\pi \in \Pi(\mathbb{P}, \mathbb{Q})} \int_{\mathbb{R}^d \times \mathbb{R}^d} d(\xi, \eta)^p \pi(d\xi, d\eta) \right)^{1/p}$$

where π is a coupling between the marginal distributions $\xi \sim \mathbb{P}$ and $\eta \sim \mathbb{Q}$, and d is a pseudometric defined on \mathcal{Z} .

Let's bound the difference between expectations under \mathbb{P} and \mathbb{P}° . For any coupling π between \mathbb{P} and \mathbb{P}° :

$$\begin{aligned}\mathbb{E}_{z \sim \mathbb{P}}[\ell(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ}[\ell(z; \theta)] &= \int \ell(\xi; \theta) \pi(d\xi, d\eta) - \int \ell(\eta; \theta) \pi(d\xi, d\eta) \\ &= \int (\ell(\xi; \theta) - \ell(\eta; \theta)) \pi(d\xi, d\eta)\end{aligned}$$

This equality holds due to the marginal properties of a coupling. Now, assume that the loss function $\ell(z; \theta)$ is $L_{\ell,z}$ -Lipschitz with respect to z (i.e., $|\ell(\xi; \theta) - \ell(\eta; \theta)| \leq L_{\ell,z} d(\xi, \eta)$ for some constant $L_{\ell,z}$). Then, for $p \geq 1$, we can apply Hölder's inequality (or the generalized mean inequality for distances):

$$\begin{aligned} |\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell(z; \theta)]| &\leq \int |\ell(\xi; \theta) - \ell(\eta; \theta)| \pi(d\xi, d\eta) \\ &\leq L_{\ell,z} \int d(\xi, \eta) \pi(d\xi, d\eta) \\ &\leq L_{\ell,z} \left(\int d(\xi, \eta)^p \pi(d\xi, d\eta) \right)^{1/p} \end{aligned}$$

Taking the supremum over all couplings π , and then over all $\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \mathcal{W}_p)$:

$$\begin{aligned} |\mathcal{L}^{\mathcal{W}_p}(\theta) - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)]| &= \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \mathcal{W}_p)} |\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)]| \\ &\leq L_{\ell,z} \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \mathcal{W}_p)} \mathcal{W}_p(\mathbb{P}, \mathbb{P}^\circ) \\ &\leq L_{\ell,z} \varepsilon_n \end{aligned}$$

Thus, we can set $\Delta_n = L_{\ell,z} \varepsilon_n$. If we choose $\varepsilon_n \asymp n^{-1}$, then $\Delta_n = O(n^{-1})$, which aligns with the condition for the $O(n^{-1/2})$ rate in the Master Theorem. \square

H Proof of "Fast Rate" KL-DRO-REBEL

Proof of Corollary 2. First, recall Lemma 9 (Gibbs variational principle characterization of the KL divergence) for probability measures \mathbb{P}, \mathbb{Q} :

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) = \sup_{g: \mathcal{Z} \rightarrow \mathbb{R}} \{\mathbb{E}_{\mathbb{P}}[g] - \log \mathbb{E}_{\mathbb{Q}}[e^g]\}$$

Let $f = \ell(z; \theta)$. For any $\lambda \geq 0$, we can choose $g(z) = \lambda(f(z) - \mathbb{E}_{\mathbb{Q}}[f])$ to obtain a lower bound for $D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q})$:

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) \geq \sup_{\lambda \geq 0} \left\{ \lambda (\mathbb{E}_{\mathbb{P}}[f] - \mathbb{E}_{\mathbb{Q}}[f]) - \log \mathbb{E}_{\mathbb{Q}} \left[e^{\lambda(f - \mathbb{E}_{\mathbb{Q}}[f])} \right] \right\}$$

Now, suppose that $f = \ell(z; \theta) \in [0, K_\ell]$ almost surely (from Lemma 12). By Hoeffding's Lemma (Lemma 10), if f is bounded in $[0, K_\ell]$, then $f - \mathbb{E}_{\mathbb{Q}}[f]$ is sub-Gaussian with parameter $K_\ell/2$. Thus, we have that

$$\log \mathbb{E}_{\mathbb{Q}} \left[e^{\lambda(f - \mathbb{E}_{\mathbb{Q}}[f])} \right] \leq \frac{\lambda^2 (K_\ell/2)^2}{2} = \frac{\lambda^2 K_\ell^2}{8}$$

Substituting this bound into the KL inequality:

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) \geq \sup_{\lambda \geq 0} \left\{ \lambda (\mathbb{E}_{\mathbb{P}}[f] - \mathbb{E}_{\mathbb{Q}}[f]) - \frac{\lambda^2 K_\ell^2}{8} \right\}$$

The expression in the curly brackets is a concave quadratic in λ . Its supremum is attained at $\lambda^* = \frac{4(\mathbb{E}_{\mathbb{P}}[f] - \mathbb{E}_{\mathbb{Q}}[f])}{K_\ell^2}$. Plugging this optimal λ^* back into the expression, we find that

$$D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q}) \geq \frac{2(\mathbb{E}_{\mathbb{P}}[f] - \mathbb{E}_{\mathbb{Q}}[f])^2}{K_\ell^2}$$

Rearranging terms, we obtain a bound on the absolute difference in expectations:

$$|\mathbb{E}_{\mathbb{P}}[f] - \mathbb{E}_{\mathbb{Q}}[f]| \leq \frac{1}{\sqrt{2}} K_\ell \sqrt{D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q})}$$

Now, taking $\mathbb{Q} = \mathbb{P}^\circ$ and $f = \ell(z; \theta)$, and considering the supremum within the KL ball $\mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \text{KL})$:

$$\begin{aligned} |\mathcal{L}^{\text{KL}}(\theta) - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)]| &= \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \text{KL})} |\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)]| \\ &\leq \frac{1}{\sqrt{2}} K_\ell \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \text{KL})} \sqrt{D_{\text{KL}}(\mathbb{P} \parallel \mathbb{P}^\circ)} \\ &\leq \frac{1}{\sqrt{2}} K_\ell \sqrt{\varepsilon_n} \end{aligned}$$

Thus, we can set $\Delta_n = \frac{1}{\sqrt{2}} K_\ell \sqrt{\varepsilon_n}$. If we choose $\varepsilon_n \asymp n^{-2}$, then $\Delta_n = O(n^{-1})$, which satisfies the condition for the $O(n^{-1/2})$ rate in the Master Theorem. \square

I Proof of "Fast Rate" χ^2 -DRO-REBEL

Proof of Corollary 3. By Theorem 8 (Hammersley-Chapman-Robbins (HCR) lower bound), we immediately have:

$$|\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \mathbb{E}_{z \sim \mathbb{P}^\circ} [\ell(z; \theta)]| \leq \sqrt{\text{Var}_{\mathbb{P}^\circ} (\ell(z; \theta)) \chi^2(\mathbb{P} \parallel \mathbb{P}^\circ)}$$

Since $\ell(z; \theta) \in [0, K_\ell]$ almost surely, its variance is bounded by $\text{Var}_{\mathbb{P}^\circ} (\ell(z; \theta)) \leq K_\ell^2/4$. Thus, we have:

$$\begin{aligned} \left| \mathcal{L}^{\chi^2}(\theta) - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)] \right| &= \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \chi^2)} |\mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \mathbb{E}_{\mathbb{P}^\circ} [\ell(z; \theta)]| \\ &\leq \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \chi^2)} \sqrt{\frac{K_\ell^2}{4} \chi^2(\mathbb{P} \parallel \mathbb{P}^\circ)} \\ &\leq \frac{K_\ell}{2} \sup_{\mathbb{P} \in \mathcal{B}_{\varepsilon_n}(\mathbb{P}^\circ; \chi^2)} \sqrt{\chi^2(\mathbb{P} \parallel \mathbb{P}^\circ)} \\ &\leq \frac{K_\ell}{2} \sqrt{\varepsilon_n} \end{aligned}$$

Thus, we can set $\Delta_n = \frac{K_\ell}{2} \sqrt{\varepsilon_n}$. If we choose $\varepsilon_n \asymp n^{-2}$, then $\Delta_n = O(n^{-1})$, which satisfies the condition for the $O(n^{-1/2})$ rate in the Master Theorem. \square

J Proof of "Fast Rate" WDPO

Proof. From the "Master Theorem", all we need to do is verify that the Wasserstein DPO objective and the DPO loss function satisfy the four conditions.

1. Verification of Local Strong Convexity From Appendix B.3, Lemma 11 of Xu et al. [2025], we know that the Wasserstein DPO loss, $\mathcal{L}^W(\theta)$, is $\gamma\lambda$ -strongly convex with respect to the Euclidean norm $\|\cdot\|_2$. This directly satisfies the first condition with a strong convexity parameter $\alpha = \gamma\lambda$ where $\gamma = \frac{\beta^2 e^{4\beta B}}{(1+e^{4\beta B})^2}$ and λ is from the data coverage assumption.

2. Verification of Lipschitz Loss (in θ) We show that the pointwise DPO loss, $\ell(z; \theta) = -y \log \sigma(\beta h_\theta) - (1-y) \log \sigma(-\beta h_\theta)$, is Lipschitz in θ . The gradient with respect to θ is $\nabla_\theta \ell(z; \theta) = \partial \ell / \partial h_\theta \cdot \nabla_\theta h_\theta$.

First, we bound the norm of the gradient of the preference score. Using the log-linear policy assumption:

$$\begin{aligned} h_\theta(s, a^1, a^2) &:= \left(\log \frac{\pi_\theta(a^1|s)}{\pi_{\text{ref}}(a^1|s)} \right) - \left(\log \frac{\pi_\theta(a^2|s)}{\pi_{\text{ref}}(a^2|s)} \right) \\ &= (\log \pi_\theta(a^1|s) - \log \pi_{\text{ref}}(a^1|s)) - (\log \pi_\theta(a^2|s) - \log \pi_{\text{ref}}(a^2|s)) \\ &= (\langle \theta, \psi(s, a^1) \rangle - \langle \theta_{\text{ref}}, \psi(s, a^1) \rangle) - (\langle \theta, \psi(s, a^2) \rangle - \langle \theta_{\text{ref}}, \psi(s, a^2) \rangle) \\ &= \langle \theta - \theta_{\text{ref}}, \psi(s, a^1) - \psi(s, a^2) \rangle \end{aligned}$$

The gradient of h_θ with respect to θ is $\nabla_\theta h_\theta = \psi(s, a^1) - \psi(s, a^2)$. Its norm is bounded:

$$\begin{aligned} \|\nabla_\theta h_\theta\|_2 &= \|\psi(s, a^1) - \psi(s, a^2)\|_2 \leq \|\psi(s, a^1)\|_2 + \|\psi(s, a^2)\|_2 \\ &\leq 2 \end{aligned}$$

Second, we bound the magnitude of the derivative of the logistic loss with respect to h_θ .

$$\frac{\partial \ell}{\partial h_\theta} = -y\beta(1 - \sigma(\beta h_\theta)) + (1-y)\beta\sigma(\beta h_\theta) = \beta((1-y)\sigma(\beta h_\theta) - y\sigma(-\beta h_\theta))$$

Since $y \in \{0, 1\}$ and $\sigma(\cdot) \in (0, 1)$, the magnitude is maximized when either term is active, giving $|\partial \ell / \partial h_\theta| \leq \beta$. Combining these results, the norm of the gradient is bounded:

$$\|\nabla_\theta \ell(z; \theta)\|_2 = \left| \frac{\partial \ell}{\partial h_\theta} \right| \cdot \|\nabla_\theta h_\theta\|_2 \leq 2\beta$$

Thus, the pointwise DPO loss is L_g -Lipschitz in θ , with $L_g = 2\beta$.

3. Verification of Bounded Loss From the derivation in Step 2 and the Cauchy-Schwarz inequality, we have

$$\begin{aligned} |h_\theta| &= |\langle \theta - \theta_{\text{ref}}, \psi(s, a^1) - \psi(s, a^2) \rangle| \\ &\leq \|\theta - \theta_{\text{ref}}\|_2 \|\psi(s, a^1) - \psi(s, a^2)\|_2 \\ &\leq 4B \end{aligned}$$

This implies the argument βh_θ is in $[-4\beta B, 4\beta B]$, and the loss is uniformly bounded by:

$$K_\ell = \log(1 + e^{4\beta B})$$

All four conditions of the Master Theorem have been verified for the Wasserstein DPO problem. We can now substitute the derived constants $\alpha = \gamma\lambda$, $L_g = 2\beta$, $K_\ell = \log(1 + e^{4\beta B})$, and $\Delta_n = n^{-1}L_{\ell,z}$ (from Appendix G) into the theorem's final bound. This yields:

$$\|\hat{\theta}_n^{\mathcal{W}_p} - \theta^{\mathcal{W}_p}\|_2 \leq \frac{32c_0\beta}{\gamma\lambda} \sqrt{\frac{d}{n}} + \frac{16c_1\beta}{\gamma\lambda} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8L_{\ell,z}}{\gamma\lambda n}} + \sqrt{\frac{4c_2 \log(1 + e^{4\beta B}) \log(1/\delta)}{\gamma\lambda n}}$$

□

K Proof of "Fast Rate" KLDPO

Proof. The proof proceeds by verifying that the KL-DPO objective satisfies the four conditions of the Master Theorem.

As we did for WDPO, we verify that the KL DPO objective satisfy the four conditions. Since we already proved that ℓ_{DPO} is Lipschitz in θ and uniformly bounded by $K_\ell = \log(1 + e^{4\beta B})$, we must just verify local strong convexity. From Appendix C, Lemma 14 of Xu et al. [2025], the KL-DPO loss, $\mathcal{L}^{KL}(\theta)$, is $\gamma\lambda$ -strongly convex with respect to the Euclidean norm $\|\cdot\|_2$. This directly satisfies the first condition with a strong convexity parameter $\alpha = \gamma\lambda$ where $\gamma = \frac{\beta^2 e^{4\beta B}}{(1 + e^{4\beta B})^2}$ and λ is from the data coverage assumption.

Thus all four conditions of the Master Theorem have been verified for the KL-DPO problem. We can now substitute the derived constants $\alpha = \gamma\lambda$, $L_g = 2\beta$, $K_\ell = \log(1 + e^{4\beta B})$, and $\Delta_n = 2^{-1/2}n^{-1}K_\ell\sqrt{\varepsilon_n}$ (from Appendix 2) into the theorem's final bound. This yields:

$$\|\hat{\theta}_n^{\text{KL}} - \theta^{\text{KL}}\|_2 \leq \frac{32c_0\beta}{\gamma\lambda} \sqrt{\frac{d}{n}} + \frac{16c_1\beta}{\gamma\lambda} \sqrt{\frac{\log(1/\delta)}{n}} + \sqrt{\frac{8 \log(1 + e^{4\beta B})}{\gamma\lambda n}} + \sqrt{\frac{4c_2 \log(1 + e^{4\beta B}) \log(1/\delta)}{\gamma\lambda n}}$$

□

L Proof of Tractable χ^2 -DRO-REBEL

Proof. Let \mathbb{P}_n be the empirical distribution. The robust optimization problem is given by:

$$\mathcal{L}_n^{\chi^2}(\theta; \rho) = \sup_{\mathbb{P}} \mathbb{E}_{\mathbb{P}}[\ell(z; \theta)] \quad \text{s.t.} \quad D_{\chi^2}(\mathbb{P} \parallel \mathbb{P}_n) \leq \rho, \quad \mathbb{P} \geq 0, \quad \mathbb{E}_{\mathbb{P}}[1] = 1.$$

The χ^2 -divergence is defined by $f(t) = \frac{1}{2}(t - 1)^2$. The Fenchel conjugate $f^*(s) = \sup_{t \geq 0} \{st - f(t)\}$. For $s \in \mathbb{R}$, $f'(t) = t - 1$. Setting $s = t - 1$, we get $t = s + 1$. Substituting this into the definition of $f^*(s)$: $f^*(s) = s(s + 1) - \frac{1}{2}((s + 1) - 1)^2 = s^2 + s - \frac{1}{2}s^2 = \frac{1}{2}s^2 + s$. This derivation holds for $t \geq 0$, which implies $s + 1 \geq 0 \implies s \geq -1$. If $s < -1$, the optimal t would be negative, violating $t \geq 0$. In this case, $f^*(s)$ becomes ∞ due to the constraint $t \geq 0$. This is implicitly handled by the dual formulation. According to Lemma 7 [Duchi and Namkoong, 2020], the dual form of the f -divergence based DRO problem is:

$$\sup_{\mathbb{P}: D_f(\mathbb{P} \parallel \mathbb{P}_n) \leq \rho} \mathbb{E}_{\mathbb{P}}[\ell(z; \theta)] = \inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \left\{ \lambda \mathbb{E}_{\mathbb{P}_n} \left[f^* \left(\frac{\ell(z; \theta) - \eta}{\lambda} \right) \right] + \lambda \rho + \eta \right\}.$$

Substituting $f^*(s) = \frac{1}{2}s^2 + s$ into this dual formulation, with $s = \frac{\ell_i - \eta}{\lambda}$:

$$\mathcal{L}_n^{\chi^2}(\theta; \rho) = \inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \left\{ \lambda \rho + \eta + \mathbb{E}_{\mathbb{P}_n} \left[\lambda \left(\frac{1}{2} \left(\frac{\ell_i - \eta}{\lambda} \right)^2 + \frac{\ell_i - \eta}{\lambda} \right) \right] \right\}.$$

This simplifies to:

$$\mathcal{L}_n^{\chi^2}(\theta; \rho) = \inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \left\{ \lambda \rho + \eta + \mathbb{E}_{\mathbb{P}_n} \left[\frac{(\ell_i - \eta)^2}{2\lambda} + (\ell_i - \eta) \right] \right\}.$$

Let $X_i = \ell_i - \eta$. The objective becomes:

$$\inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \left\{ \lambda \rho + \eta + \frac{1}{n} \sum_{i=1}^n \left[\frac{X_i^2}{2\lambda} + X_i \right] \right\}.$$

The non-negativity constraint $P(z_i) \geq 0$ in the primal problem implies $w_i \geq 0$, which means $1 + \frac{\ell_i - \eta}{\lambda} \geq 0$. This is equivalent to $\frac{\ell_i - \eta}{\lambda} \geq -1$. This constraint is handled by a special form of f^* or by considering the dual's objective piece-wise. When $1 + \frac{\ell_i - \eta}{\lambda} < 0$, the optimal $w_i = 0$, which essentially means that this instance z_i is excluded from the worst-case distribution. This leads to the emergence of the positive part $(\cdot)_+$ in the objective. Specifically, for χ^2 -divergence, it is a known result in robust optimization that the problem is equivalent to:

$$\mathcal{L}_n^{\chi^2}(\theta; \rho) = \inf_{\eta \in \mathbb{R}} \left\{ \eta + \inf_{\lambda > 0} \left\{ \lambda \rho + \frac{1}{n} \sum_{i=1}^n \frac{(\ell_i - \eta)_+^2}{2\lambda} \right\} \right\}.$$

The term $\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)]$ or $\frac{1}{n} \sum X_i$ from the initial general dual disappears because when minimizing over λ , the optimal choice of λ ensures that these terms effectively vanish or are absorbed into the structure involving the positive part, reflecting the boundary conditions where $w_i = 0$. This transformation from the general f -divergence dual to the specific one-dimensional \inf_{η} form for χ^2 -divergence (with the $(\cdot)_+^2$ term) is detailed in Namkoong and Duchi [2017a] (Lemma 2.1 and Appendix A.1).

Now, we solve the inner minimization with respect to λ for a fixed η . Let $Y_i = (\ell_i - \eta)_+$. The inner objective is:

$$G(\lambda) = \lambda \rho + \frac{1}{n} \sum_{i=1}^n \frac{Y_i^2}{2\lambda}.$$

To find the optimal λ^* , we differentiate $G(\lambda)$ with respect to λ and set it to zero:

$$\frac{dG(\lambda)}{d\lambda} = \rho - \frac{1}{n} \sum_{i=1}^n \frac{Y_i^2}{2\lambda^2} = 0.$$

Solving for λ^2 :

$$\lambda^2 = \frac{\sum_{i=1}^n Y_i^2}{2n\rho} = \frac{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}.$$

Since $\lambda > 0$ and $\rho > 0$, we take the positive square root:

$$\lambda^* = \sqrt{\frac{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}}.$$

Substitute λ^* back into the inner objective $G(\lambda)$:

$$\begin{aligned}
G(\lambda^*) &= \sqrt{\frac{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}} \cdot \rho + \frac{1}{n} \sum_{i=1}^n \frac{(\ell_i - \eta)_+^2}{2\sqrt{\frac{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}}} \\
&= \rho \sqrt{\frac{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}} + \frac{1}{2} \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2] \sqrt{\frac{2\rho}{\mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}} \\
&= \sqrt{\frac{\rho^2 \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2\rho}} + \frac{1}{2} \sqrt{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]} \\
&= \sqrt{\frac{\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{2}} + \frac{1}{2} \sqrt{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]} \\
&= \sqrt{\frac{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{4}} + \sqrt{\frac{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{4}} \\
&= 2\sqrt{\frac{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}{4}} \\
&= \sqrt{2\rho \mathbb{E}_{\mathbb{P}_n}[(\ell_i - \eta)_+^2]}.
\end{aligned}$$

Therefore, the robust objective simplifies to:

$$\mathcal{L}_n^{\chi^2}(\theta; \rho) = \inf_{\eta \in \mathbb{R}} \left\{ \eta + \sqrt{\frac{2\rho}{n} \sum_{i=1}^n (\ell_i - \eta)_+^2} \right\},$$

which matches the first part of the proposition. For the worst-case weights, the optimal density ratio w_i is found by $w_i = \max\left(0, 1 + \frac{\ell_i - \eta^*}{\lambda^*}\right)$. This can be rewritten using the specific structure of χ^2 -divergence as $w_i = \frac{(\ell_i - \eta^*)_+}{n \lambda^*}$ where λ^* is the optimal λ for the given η^* . This form ensures that $\sum w_i = n$ (as demonstrated in Namkoong and Duchi [2017a], Appendix A.1) and aligns with the intuitive notion that higher loss values contribute more to the worst-case distribution.

We now prove that this problem can be solved efficiently by first establishing the convexity of

$$f(\eta) := \eta + \sqrt{\frac{2\varepsilon_n}{n} \sum_{i=1}^n (\ell_i - \eta)_+^2}$$

and show that the search space for its minimum is bounded.

1. Convexity of the Objective Function We show that $f(\eta)$ is a convex function of η . The function can be written as the sum of two functions, $f(\eta) = g(\eta) + h(\eta)$, where $g(\eta) = \eta$ and $h(\eta) = \sqrt{C \sum_{i=1}^n v_i(\eta)^2}$ with $C = \frac{2\varepsilon_n}{n}$ and $v_i(\eta) = (\ell_i - \eta)_+$.

- The function $g(\eta) = \eta$ is linear and therefore convex.
- For each i , the function $v_i(\eta) = \max(0, \ell_i - \eta)$ is a hinge function, which is the maximum of two affine (and thus convex) functions, 0 and $\ell_i - \eta$. Therefore, each $v_i(\eta)$ is convex in η .
- Let $v(\eta) = [v_1(\eta), \dots, v_n(\eta)]^\top$ be a vector-valued function. Since each component is convex, the function $v(\eta)$ is convex.
- The function $\phi(v) = \sqrt{C} \|v\|_2$ is the scaled L2-norm, which is a convex function. Furthermore, for the domain of v (which has non-negative components), ϕ is a non-decreasing function with respect to each component.
- The function $h(\eta) = \phi(v(\eta))$ is a composition of a convex function ϕ and a convex vector-valued function $v(\eta)$. A composition of this form preserves convexity.
- Since $f(\eta)$ is the sum of two convex functions, $g(\eta)$ and $h(\eta)$, it is itself a convex function.

2. Characterization of the Optimum Since $f(\eta)$ is convex, a point η^* is a minimum if and only if the subgradient contains zero, i.e., $0 \in \partial f(\eta^*)$. The subdifferential of f is $\partial f(\eta) = 1 + \partial h(\eta)$. For any point $\eta \in \mathbb{R}$ where f is differentiable (i.e., $\eta \neq \ell_i$ for all i where $\ell_i > \eta$), the derivative is given by:

$$f'(\eta) = 1 - \sqrt{\frac{2\varepsilon_n}{n}} \cdot \frac{\sum_{i:\ell_i > \eta} (\ell_i - \eta)}{\sqrt{\sum_{i:\ell_i > \eta} (\ell_i - \eta)^2}}$$

Because $f(\eta)$ is convex, its subgradient is a monotonically non-decreasing function of η .

3. Bounded Search Space We can establish a finite upper bound for the search space. Consider any $\eta > \max_i \{\ell_i\}$. For such an η , the term $(\ell_i - \eta)_+ = 0$ for all $i = 1, \dots, n$. The objective function simplifies to:

$$f(\eta) = \eta, \quad \text{for } \eta > \max_i \{\ell_i\}$$

In this region, the derivative is $f'(\eta) = 1$. Since the function is strictly increasing for all $\eta > \max_i \{\ell_i\}$, the minimizer η^* must satisfy:

$$\eta^* \leq \max_i \{\ell_i\}$$

This provides a concrete upper bound for the search. A lower bound can also be established, as $f(\eta) \rightarrow \infty$ when $\eta \rightarrow -\infty$. Thus, the search for the minimum can be restricted to a finite interval.

4. Algorithm and Complexity The properties above guarantee that we can find the unique minimizer η^* efficiently. The monotonicity of the subgradient allows the use of a binary search algorithm.

1. Define a search interval $[L, U]$, where $U = \max_i \{\ell_i\}$ and L is a sufficiently small lower bound.
2. At each iteration, select a candidate $\eta_c = (L + U)/2$.
3. Compute a subgradient $g_c \in \partial f(\eta_c)$. This takes $O(n)$ time as it requires summing over the n loss terms.
4. If $g_c > 0$, the minimum must lie to the left, so we set $U = \eta_c$.
5. If $g_c < 0$, the minimum must lie to the right, so we set $L = \eta_c$.

This procedure is repeated until the interval $[L, U]$ is sufficiently small. The number of iterations required to achieve a desired precision ϵ is $O(\log((U - L)/\epsilon))$. The total complexity of this search is $O(n \cdot \log(1/\epsilon))$. For Algorithm 4, if we assume that $\text{Card}(\{\ell_i\}_{i=1}^n) = n$, then the runtime will be $O(n \log n)$. \square