

THE PERILS OF OPTIMIZING LEARNED REWARD FUNCTIONS: LOW TRAINING ERROR DOES NOT GUARANTEE LOW REGRET

Anonymous authors

Paper under double-blind review

For reviewers:

Red parts are from the original submission and will be removed for the camera-ready version.

Blue parts are new additions from the rebuttal phase

ABSTRACT

In reinforcement learning, specifying reward functions that capture the intended task can be very challenging. Reward learning aims to address this issue by *learning* the reward function. However, a learned reward model may have a low error on the data distribution, and yet subsequently produce a policy with large regret. We say that such a reward model has an *error-regret mismatch*. The main source of an error-regret mismatch is the distributional shift that commonly occurs during policy optimization. In this paper, we mathematically show that a sufficiently low expected test error of the reward model guarantees low worst-case regret, but that for any *fixed* expected test error, there exist realistic data distributions that allow for error-regret mismatch to occur. We then show that similar problems persist even when using policy regularization techniques, commonly employed in methods such as RLHF. **Our theoretical results highlight the importance of developing new ways to measure the quality of learned reward models. We hope our results stimulate the theoretical and empirical study of improved methods to learn reward models, and better ways to reliably measure their quality.**

1 INTRODUCTION

To solve a sequential decision problem with reinforcement learning (RL), we must first formalize that decision problem using a *reward function* (Sutton & Barto, 2018). However, for complex tasks, reward functions are often hard to specify correctly. To solve this problem, it is increasingly popular to *learn* reward functions with *reward learning algorithms*, instead of specifying the reward functions manually. There are many different reward learning algorithms (e.g., Ng & Russell, 2000; Tung et al., 2018; Brown & Niekum, 2019; Palan et al., 2019), with one of the most popular being *reward learning from human feedback* (RLHF) (Christiano et al., 2017; Ibarz et al., 2018).

For any learning algorithm, it is a crucial question whether or not that learning algorithm is guaranteed to converge to a “good” solution. For example, in the case of supervised learning for classification, it can be shown that a learning algorithm that produces a model with a low *empirical error* (i.e., training error) is likely to have a low *expected error* (i.e., test error), given a sufficient amount of training data and assuming that both the training data and the test data is drawn i.i.d. from a single stationary distribution (Kearns & Vazirani, 1994). In the case of normal supervised learning and standard assumptions, we can therefore be confident that a learning algorithm will converge to a good model, provided that it is given a sufficient amount of training data.

Since reward models are also typically learned by supervised learning, we might assume that classical learning-theoretic guarantees carry over. However, these guarantees only ensure that the reward model is approximately correct *relative to the training distribution*. But after reward learning, we optimize a policy to maximize the learned reward, which effectively leads to a *distributional shift*. This raises the worry that the trained policy can exploit regions of the state space with abnormally

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

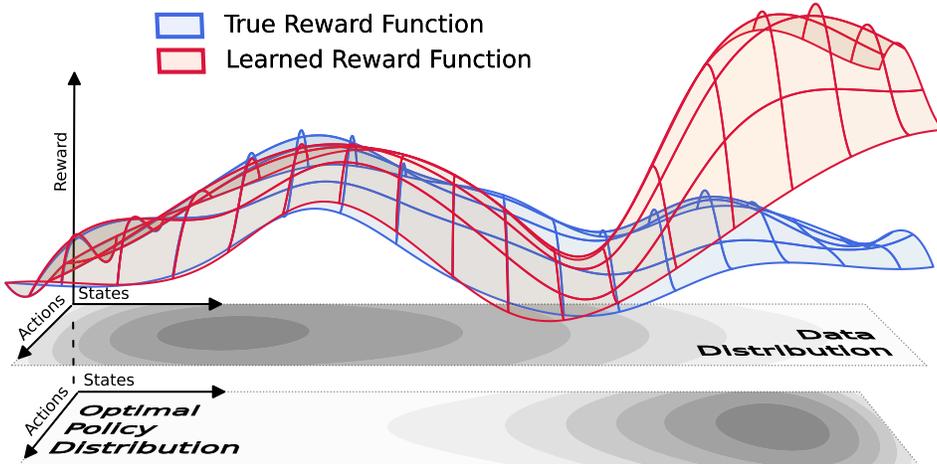


Figure 1: Reward models (red function) are commonly trained by supervised learning to approximate some latent, true reward (blue function). Given enough data, one can hope that the reward model is close to the true reward function on average over the data distribution (upper gray layer) — the expected *error* is low. However, low expected error only guarantees a good approximation to the true reward function in areas with high coverage by the data distribution! On the other hand, optimizing an RL policy to maximize the learned reward model induces a distribution shift which can lead the policy to exploit uncertainties of the learned reward model in low-probability areas of the transition space (lower gray layer). This may then lead to high *regret*. We refer to this phenomenon as *error-regret mismatch*.

high learned rewards if those regions have a low data coverage during training. In this case, we can have reward models that have both a low error on the training distribution and an optimal policy with large regret, a phenomenon we call *error-regret mismatch*. We visualize this concern in Figure 1.

To single out the issue of error-regret mismatch in our analysis, we take the goals of classical learning theory as a given and show that *they are not enough to ensure low regret*. More precisely, in probably approximately correct (PAC) learning (Kearns & Vazirani, 1994) the goal is to derive a sample size that guarantees a certain likelihood (“P”) of an approximately correct (“AC”) model on new data points sampled from the training distribution. In our results, we assume that we *already have* an approximately correct reward model on a data distribution, and then investigate what we can or can not conclude about the regret of policies trained to maximize the modeled reward.

Our theoretical analysis shows that guarantees in policy regret are very sensitive to the data distribution used to train the reward model, leading to our notions of *safe* and *unsafe data distributions*. Moreover, we find evidence that some MDPs are in a certain sense “too large” to allow for safe data distributions relative to a reasonable reward model error and desired regret bound. We establish for general MDPs:

1. As the error of a learned reward model on a data distribution goes to zero, the worst-case regret of optimizing a policy according to that reward model also goes to zero (Propositions 3.1 and 3.2)
2. However, for any $\epsilon > 0$, whenever a data distribution has sufficiently low coverage of some bad policy, it is *unsafe*; in other words, there exists a reward model that achieves an expected error of ϵ but has a high-regret optimal policy (Proposition 3.3), a case of error-regret mismatch.
3. As a consequence, when an MDP has a large number of independent bad policies, *every* data distribution is unsafe (Corollary 3.4).
4. More precisely, we derive a set of linear constraints that precisely characterize the safe data distributions for a given MDP (Theorem 3.5).

We then investigate the case of *regularized* policy optimization (including KL-regularized policy optimization, which is commonly used in methods such as RLHF). We derive regularized versions

of Propositions 3.1 and 3.3 in Proposition 4.1 and Theorem 4.2. This shows that regularization alone is no principled solution to error-regret mismatch.

We then develop several generalizations of our results for different types of data sources for reward model training, such as preferences over trajectories (Propositions 5.2 and 5.3), and trajectory scoring (Proposition 5.1). Lastly, motivated by the recent success of large language models (OpenAI, 2022; Gemini Team, 2023; Anthropic, 2023), we provide an analysis for the special case of RLHF in the contextual bandit case where we prove a stronger version (Theorem 6.1) of the failure mode already discussed in Theorem 4.2 for general MDPs.

1.1 RELATED WORK

Note: We provide a more extensive related work section in Appendix A

Reward Learning Reward learning is a key concept in reinforcement learning that involves learning the reward function for complex tasks with latent and difficult-to-specify reward functions. Many methods have been developed to incorporate various types of human feedback into the reward learning process (Wirth et al., 2017; Ng et al., 2000; Bajcsy et al., 2017; Jeon et al., 2020).

Challenges in Reward Learning Reward learning presents several challenges (Casper et al., 2023; Lang et al., 2024b; Skalse & Abate, 2023; 2024), such as *reward misgeneralization*, where a learned reward model performs well on in-distribution data but misgeneralizes on out-of-distribution data (Skalse et al., 2023). This can lead to unintended consequences in real-world applications.

Reward misgeneralization can also result in *reward hacking* (Krakovna, 2020), a consequence of Goodhart’s law (Goodhart, 1984; Zhuang & Hadfield-Menell, 2020; Hennessy & Goodhart, 2023; Strathern, 1997; Karwowski et al., 2023). Reward hacking has been extensively studied both theoretically (Skalse et al., 2022; 2024; Zhuang & Hadfield-Menell, 2020) and empirically (Zhang et al., 2018; Farebrother et al., 2018; Cobbe et al., 2019; Krakovna, 2020; Gao et al., 2023; Tien et al., 2022).

Distributional Shifts in policy learning During policy optimization, a distribution shift occurs where the policy under training can explore areas of the input space that are outside of the reward model’s training distribution. This might lead to large regret in case that the reward model is misgeneralized. To address issues with distributional shifts in reward learning specifically, prior work has proposed many different methods, such as ensembles of conservative reward models (Coste et al., 2023), averaging weights of multiple reward models (Ramé et al., 2024), iteratively updating training labels (Zhu et al., 2024), on-policy reward learning (Lang et al., 2024a), and distributionally robust planning (Zhan et al., 2023).

Worst-case regret studies Several prior works perform theoretical investigations of policy regret under a worst-case MDP in settings involving imitation learning (Ross et al., 2011), offline RL (Kim et al., 2024; Jin et al., 2021; Zhang et al., 2022), and other RL settings (Lu et al., 2024; Laidlaw et al., 2024; Kwa et al., 2024). Furthermore, additional work performs analyses of RLHF (Cen et al., 2024; Xiong et al., 2024; Zhu et al., 2023; Ji et al., 2023; Mehta et al., 2023) and reward learning in general (Agarwal et al., 2012; Foster et al., 2020) in the contextual bandit setting.

Our work is most closely related to (Zhu et al., 2024; Nika et al., 2024; Cen et al., 2024), which provide examples of high regret and regret bounds specifically for different RLHF settings. We contrast ourselves from their work by focusing on theoretically analyzing the guarantees from errors with respect to data distributions on the regret of the final policy. We thus narrow down and illuminate issues that are unique to RL with policy learning compared to supervised learning. In doing so, we assume *arbitrary* MDPs and investigate the regret a policy might attain for a worst-case reward model (instead of a worst-case MDP).

2 PRELIMINARIES

A *Markov Decision Process* (MDP) is a tuple $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ where \mathcal{S} is a set of *states*, \mathcal{A} is a set of *actions*, $\tau : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is a *transition function*, $\mu_0 \in \Delta(\mathcal{S})$ is an *initial state distribution*,

162 $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a *reward function*, and $\gamma \in (0, 1)$ is a *discount rate*. We define the *range* of a
 163 reward function R as $\text{range } R := \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} R(s, a) - \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} R(s, a)$.

164
 165 A *policy* is a function $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$. We denote the set of all policies by Π . A *trajectory*
 166 $\xi = \langle s_0, a_0, s_1, a_1, \dots \rangle$ is a possible path in an MDP. The *return function* G gives the cumulative
 167 discounted reward of a trajectory, $G(\xi) = \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)$, and the *evaluation function* J gives the
 168 expected trajectory return given a policy, $J(\pi) = \mathbb{E}_{\xi \sim \pi} [G(\xi)]$. A policy maximizing J is an *optimal*
 169 *policy*. We define the *regret* of a policy π with respect to reward function R as

$$170 \text{Reg}^R(\pi) := \frac{\max_{\pi' \in \Pi} J_R(\pi') - J_R(\pi)}{\max_{\pi' \in \Pi} J_R(\pi') - \min_{\pi' \in \Pi} J_R(\pi')} \in [0, 1].$$

171
 172 Here, J_R is the policy evaluation function for R .

173
 174 In this paper, we assume that \mathcal{S} and \mathcal{A} are finite, and that all states are reachable under τ and μ_0 . We
 175 also assume that $\max J_R - \min J_R \neq 0$ (since the reward function would otherwise be trivial). Note
 176 that this implies that $\text{range } R > 0$, and that $\text{Reg}^R(\pi)$ is well-defined.

177 The *state-action occupancy measure* is a function $\eta : \Pi \rightarrow \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$ mapping each policy $\pi \in \Pi$
 178 to the corresponding "state-action occupancy measure", describing the discounted frequency that
 179 each state-action tuple is visited by a policy. Formally, $\eta(\pi)(s, a) = \eta^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t \cdot P(s_t =$
 180 $s, a_t = a \mid \xi \sim \pi)$. Note that by writing the reward function R as a vector $\vec{R} \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$, we can
 181 split J into a function that is linear in R : $J(\pi) = \eta^\pi \cdot \vec{R}$. By normalizing a state-action occupancy
 182 measure η^π we obtain a *policy-induced distribution* $D^\pi := (1 - \gamma) \cdot \eta^\pi$.

183 2.1 PROBLEM FORMALIZATION OF RL WITH REWARD LEARNING

184
 185 In RL with reward learning, we assume that we have an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ where the reward
 186 function R is unknown. We may also assume that τ and μ_0 are unknown, as long as we can sample
 187 from them (though \mathcal{S} , \mathcal{A} , and γ must generally be known, at least implicitly). We then first learn a
 188 reward model \hat{R} that approximates the true reward R and then optimize a policy $\hat{\pi}$ to maximize \hat{R} .
 189 The aim of this two-step procedure is for $\hat{\pi}$ to achieve low regret under the true reward function R .
 190 We now formalize these aspects in detail for our theoretical analysis, with a visualization provided in
 191 Figure 2:

192 **Reward learning** We first learn a reward model \hat{R} from data. There are many possible data sources
 193 for reward learning, like demonstrations (Ng & Russell, 2000), preferences over trajectories (Chris-
 194 tiano et al., 2017), or even the initial environment state (Shah et al., 2019); a taxonomy can be found
 195 in (Jeon et al., 2020). Since we are concerned with problems that remain even when the reward
 196 model is already *approximately correct*, we abstract away the data sources and training procedures
 197 and assume that we learn a reward model \hat{R} which satisfies

$$198 \mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s, a) - R(s, a)|}{\text{range } R} \right] \leq \epsilon \quad (1)$$

199
 200 for some $\epsilon > 0$ and stationary distribution D over transitions $\mathcal{S} \times \mathcal{A}$. Note that this is the true
 201 expectation under D , rather than an estimate of this expectation based on some finite sample. We
 202 divide by $\text{range } R$, since the absolute error ϵ is only meaningful relative to the overall scale of the
 203 reward R .
 204

205 To be clear, most reward learning algorithms *cannot guarantee* a bound as in Equation (1) since most
 206 realistic data sources do not determine the true reward function, even for infinite data (Skalse et al.,
 207 2023). Instead, we choose Equation (1) because it serves as an *upper bound* to many common reward
 208 learning training objectives (see Appendix C.5). Thus, when we show in later sections that high regret
 209 is possible even when this inequality holds, then this problem can be expected to generalize to other
 210 data sources. We make this generalization precise for some data sources in Section 5. In particular,
 211 we will show that Equation (1) implies a low cross-entropy error between the choice distributions of
 212 the true reward function and the reward model, as is commonly used for RLHF, e.g. in the context of
 213 language models (Ziegler et al., 2019).
 214

215 **Policy optimization** Given \hat{R} , we then learn a policy $\hat{\pi}$ by solving the MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, \hat{R}, \gamma \rangle$.
 In the most straightforward case, we do this by simply finding a policy that is optimal according

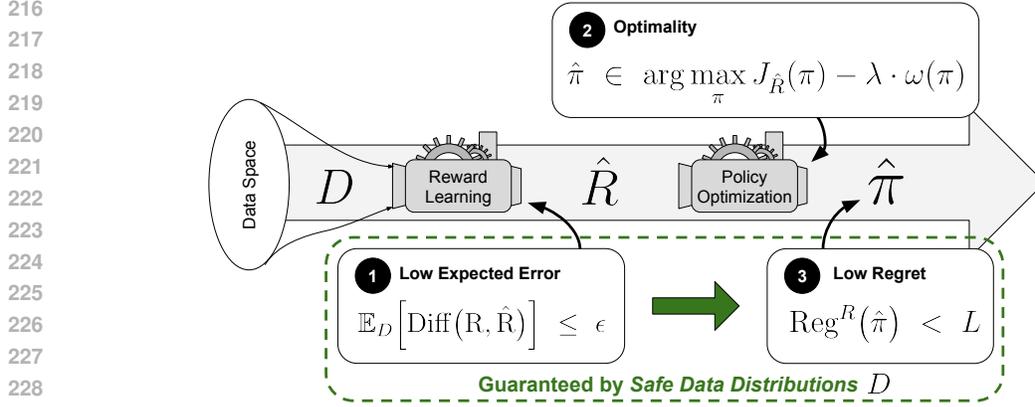


Figure 2: An abstract model of the classical reward learning pipeline. A reward model \hat{R} is trained to approximate the true reward function R under some data distribution D . The training process converges when \hat{R} is similar to R in expectation (see 1). In the second step, a policy $\hat{\pi}$ is trained to achieve high learned reward, possibly involving a regularization (see 2). We are interested in the question of when exactly this training process guarantees that $\hat{\pi}$ has low regret. More formally, we call a data distribution D *safe* whenever the implication 1 \implies 3 holds for all reward models \hat{R} that satisfy 1.

to \hat{R} . However, it is also common to perform *regularized optimization*. In that case, we make use of an additional regularization function $\omega : \Pi \rightarrow \mathbb{R}$, with $\omega(\pi) \geq 0$ for all $\pi \in \Pi$. Given \hat{R} , a regularization function ω , and a regularization weight $\lambda \in [0, \infty)$, we say that $\hat{\pi}$ is (λ, ω) -optimal if

$$\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \omega(\pi). \quad (2)$$

Typically, λ punishes large deviations from some reference policy π_{ref} , e.g. with the regularization function given by the KL-divergence $\omega(\pi) = \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}})$. π_{ref} may also be used to collect training data for the reward learning algorithm, in which case we may assume $D = D^{\pi}$ in Equation (1). Most of our results do not depend on these specific instantiations, however.

Regret minimization The aim of the previous two steps is for the policy $\hat{\pi}$ to have low regret $\text{Reg}^R(\hat{\pi})$ under the true reward function R . Our question is thus if and when it is sufficient to ensure that \hat{R} satisfies Equation 1, in order to guarantee that a policy $\hat{\pi}$ optimal according to Equation (2) has low regret $\text{Reg}^R(\hat{\pi})$.

2.2 SAFE DATA DISTRIBUTIONS

We now make the elaborations from the previous subsections more concrete by providing a formal definition of a *safe data distribution*. In particular, we say that a data distribution D is *safe*, whenever it holds that for every reward model \hat{R} that satisfies Equation (1) for D , all optimal policies of \hat{R} have low regret. We provide a visualization of this concept in Figure 2 and a formal definition in Definition 2.1.

Definition 2.1 (Safe- and unsafe data distributions). For a given MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, let $\epsilon > 0$, $L \in [0, 1]$, and $\lambda \in [0, \infty)$. Let ω be a continuous function with $\omega(\pi) \geq 0$ for all $\pi \in \Pi$. Then the set of *safe data distributions* $\text{safe}(R, \epsilon, L, \lambda, \omega)$ is the set of all distributions $D \in \Delta(\mathcal{S} \times \mathcal{A})$ such that for all possible reward models $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and policies $\hat{\pi} : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that satisfy the following two properties:

1. **Low expected error:** \hat{R} is ϵ -close to R under D , i.e., $\mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s,a) - R(s,a)|}{\text{range } \hat{R}} \right] \leq \epsilon$.
2. **Optimality:** $\hat{\pi}$ is (λ, ω) -optimal with respect to \hat{R} , i.e. $\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \omega(\pi)$.

we can guarantee that $\hat{\pi}$ has regret smaller than L , i.e.:

270 3. **Low regret:** $\hat{\pi}$ has a regret smaller than L with respect to R , i.e., $\text{Reg}^R(\hat{\pi}) < L$.

271 Similarly, we define the set of *unsafe data distributions* to be the complement of $\text{safe}(R, \epsilon, L, \lambda, \omega)$:

$$272 \quad \text{unsafe}(R, \epsilon, L, \lambda, \omega) := \{D \in \Delta(\mathcal{S} \times \mathcal{A}) \mid D \notin \text{safe}(R, \epsilon, L, \lambda, \omega)\}.$$

273 Thus, $\text{unsafe}(R, \epsilon, L, \lambda, \omega)$ consists of the data distributions D for which there *exists* a reward model
274 \hat{R} that is ϵ -close to R and a policy $\hat{\pi}$ that is (λ, π) -optimal with respect to \hat{R} , but such that $\hat{\pi}$ has large
275 regret $\text{Reg}^R(\hat{\pi}) \geq L$. In this sense, we are operating under a worst-case framework for the reward
276 model and policy learned by our training algorithms. Lastly, whenever we consider the unregularized
277 case ($\lambda = 0$ or $\omega = 0$), we drop the λ and ω to ease the notation and just use $\text{safe}(R, \epsilon, L)$ and
278 $\text{unsafe}(R, \epsilon, L)$ instead.

279 *Note: Throughout this paper, we will use the terminology that a data distribution D “allows for
280 error-regret mismatch” as a colloquial term to express that $D \in \text{unsafe}(R, \epsilon, L, \lambda, \omega)$.*

281 3 ERROR-REGRET MISMATCH FOR UNREGULARIZED POLICY OPTIMIZATION

282 In this section, we investigate the case where no regularization is used in the policy optimization stage.
283 We seek to determine if it is sufficient for a reward model to be close to the true reward function on a
284 data distribution in order to ensure low regret for the learned policy.

285 In our first result, we show that under certain conditions, a low expected error ϵ does indeed guarantee
286 that policy optimization will yield a policy with low regret.

287 **Proposition 3.1.** *Let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an arbitrary MDP, let $L \in (0, 1]$, and let $D \in \Delta(\mathcal{S} \times \mathcal{A})$
288 be a positive data distribution (i.e., a distribution such that $D(s, a) > 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$). Then
289 there exists an $\epsilon > 0$ such that $D \in \text{safe}(R, \epsilon, L)$.*

290 The proof of Proposition 3.1 can be found in Appendix D.1 (see Corollary D.7) and is based on an
291 application of Berge’s maximum theorem (Berge, 1963), and the fact that the expected distance
292 between the true reward function and the learned reward model under D is induced from a norm. See
293 Theorem 6.1 for a similar result in which the expected error in rewards is replaced by an expected
294 error in choice probabilities.

295 One might be inclined to conclude that the guarantee of Proposition 3.1 allows one to practically
296 achieve low regret by ensuring a low error ϵ (as measured by Equation 1). However, in the following
297 result we provide a more detailed analysis that shows that low regret requires a prohibitively low ϵ :

298 **Proposition 3.2.** *Let the setting be as in Proposition 3.1. If $\epsilon > 0$ satisfies*

$$299 \quad \epsilon < \frac{1 - \gamma}{\sqrt{2}} \cdot \frac{\text{range } J^R}{\text{range } R} \cdot \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a) \cdot L$$

300 *then $D \in \text{safe}(R, \epsilon, L)$.*

301 The proof can be found in Theorem D.11, Appendix D.2. Example D.13 shows that the bound
302 on ϵ is tight up to a factor of $\sqrt{2}$. This result is problematic in practice due to the dependence on
303 the minimum of D . Realistic MDPs usually contain a massive amount of states and actions, which
304 necessarily requires D to give a very small support to at least some transitions. The dependence of
305 the upper bound on D also shows that there is no ϵ for which every distribution D is guaranteed to be
306 safe, as $\min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a)$ can be arbitrarily small. We concretize this intuition by showing that in
307 every MDP and for every $\epsilon > 0$, there exist weak assumptions for which a data distribution allows
308 for a large error-regret mismatch.

309 **Proposition 3.3.** *Let $M = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP, $D \in \Delta(\mathcal{S} \times \mathcal{A})$ a data distribution,
310 $\epsilon > 0$, and $L \in [0, 1]$. Assume there exists a policy $\hat{\pi}$ with the property that $\text{Reg}^R(\hat{\pi}) \geq L$ and
311 $D(\text{supp } D^{\hat{\pi}}) < \epsilon$, where $\text{supp } D^{\hat{\pi}}$ is defined as the set of state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$ such that
312 $D^{\hat{\pi}}(s, a) > 0$. In other words, there is a “bad” policy for R that is not very supported by D . Then,
313 D allows for error-regret mismatch to occur, i.e., $D \in \text{unsafe}(R, \epsilon, L)$.*

314 The proof of Proposition 3.3 can be found in Appendix C.2 (see Proposition C.5). The intuition is
315 straightforward: There exists a reward model \hat{R} that is very similar to the true reward function R

outside the support of $D^{\hat{\pi}}$ but has very large rewards for the support of $D^{\hat{\pi}}$. Because $D(\text{supp } D^{\hat{\pi}})$ is very small, this still allows \hat{R} to have a very small expected error w.r.t. to D , while $\hat{\pi}$, the optimal policy for \hat{R} , will have regret at least L . To avoid confusions, we show in Proposition C.7 that the assumptions on ϵ in Proposition 3.2 and Proposition 3.3 cannot hold simultaneously. This is as expected since otherwise the *conclusions* of these propositions would imply that a data distribution can be both safe and unsafe.

Note that the conditions for unsafe data distributions in Proposition 3.3 also cover positive data distributions (that we showed to be eventually safe for small enough ϵ in Proposition 3.1). Furthermore, especially in very large MDPs, it is very likely that the data distribution will not sufficiently cover large parts of the support of some policies, especially since the number of (deterministic) policies grows exponentially with the number of states. Sometimes, this can lead to *all* data distributions being unsafe, as we show in the following corollary:

Corollary 3.4. *Let $M = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP, $\epsilon > 0$, and $L \in [0, 1]$. Assume there exists a set of policies Π_L with:*

- $\text{Reg}^R(\pi) \geq L$ for all $\pi \in \Pi_L$;
- $\text{supp } D^\pi \cap \text{supp } D^{\pi'} = \emptyset$ for all $\pi, \pi' \in \Pi_L$; and
- $|\Pi_L| \geq 1/\epsilon$.

Then $\text{unsafe}(R, \epsilon, L) = \Delta(\mathcal{S} \times \mathcal{A})$, i.e.: all distributions are unsafe.

The proof of Corollary 3.4 can be found in Appendix C.2 (see Corollary C.6).

Corollary 3.4 outlines sufficient conditions for a scenario where all possible data distributions are unsafe for a given MDP. This happens when there exist *many* different policies with large regret and disjoint support, which requires there to be a large action space. This could for example happen in the case of a language model interacting with a user. There are some ways to interact with the user that have large regret $\geq L$, e.g., by providing instructions for building weapons. Now consider that for a single such policy, we can easily imagine many adaptations that all behave in essentially the same way, but have individually consistent and mutually distinct writing styles or idiosyncratic differences in word choice. The set of all these variations Π_L will then be large ($|\Pi_L| > 1/\epsilon$) and consist of policies with high regret $\geq L$ that have mutually disjoint support since their actions, i.e. responses, consistently differ in style.

An example of such an MDP is the natural language environment, where the state space is a sequence of user-prompts and assistant responses. Therefore, while Proposition 3.1 shows that for every data distribution D there exists a test error ϵ small enough such that D becomes safe, Corollary 3.4 shows that in some settings ϵ might need to be extremely small in order to allow for safety.

We conclude by stating the main result of this section, which unifies all previous results and derives the most general conditions, i.e. *necessary and sufficient* conditions, for when exactly a data distribution allows for error-regret mismatch to occur:

Theorem 3.5. *For all MDPs $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and $L \in [0, 1]$, there exists a matrix M such that for all $\epsilon > 0$ and $D \in \Delta(\mathcal{S} \times \mathcal{A})$ we have:*

$$D \in \text{safe}(R, \epsilon, L) \iff M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}, \quad (3)$$

where we use the vector notation of D , and $\mathbf{1}$ is a vector containing all ones.

The proof of Theorem 3.5 can be found in Appendix C.3 (see Theorem C.16) and largely relies on geometric arguments that arise from comparing the set of unsafe reward models and the set of reward models that are close to the true reward function. Interestingly, this means that the set of *safe* data distributions resembles a polytope, in the sense that it is a convex set and is defined by the intersection of an open polyhedral set (defined by the system of strict inequalities $M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}$), and the closed data distribution simplex.

While Theorem 3.5 only proves the existence of such a matrix M , we provide further results and analyses in the appendix, namely:

1. In Appendix C.3.2 we derive closed-form expressions of the rows of matrix M , and show that its entries depend on multiple factors, such as the original reward function R , the state transition distribution τ , and the set of deterministic policies that achieve regret at least L .
2. In Appendix C.3.3 we provide an algorithm to compute matrix M .
3. In Appendix C.3.4 we provide a worked example of computing and visualizing the set of safe distributions for a toy example.

Lastly, we note that M does not depend on ϵ , and M only contains non-negative entries (see Appendix C.3.2). This allows us to recover Proposition 3.1, since by letting ϵ approach zero, the set of data distributions that fulfill the conditions in Equation (3) approaches the entire data distribution simplex. On the other hand, the dependence of M on the true reward function and the underlying MDP implies that computing M is infeasible in practice since many of these components are not known, restricting the use of M to theoretical analysis.

4 ERROR-REGRET MISMATCH FOR REGULARIZED POLICY OPTIMIZATION

In this section, we investigate the error-regret mismatch for regularized policy optimization. **We begin by showing that there are conditions under which a low expected error ϵ guarantees that a provided data distribution is safer than an initial reference policy:** First, we prove that for almost any reference policy π_{ref} that achieves regret L and minimizes the regularization term ω , there exists a sufficiently small ϵ such that reward learning within ϵ of the true reward function preserves the regret bound L .

Proposition 4.1. *Let $\lambda \in (0, \infty)$, let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be any MDP, and let $D \in \mathcal{S} \times \mathcal{A}$ be any data distribution that assigns positive probability to all transitions. Let $\omega : \Pi \rightarrow \mathbb{R}$ be a continuous regularization function that has a reference policy π_{ref} as a minimum.¹ Assume that π_{ref} is not (λ, ω) -optimal for R and let $L = \text{Reg}^R(\pi_{\text{ref}})$. Then there exists $\epsilon > 0$ such that $D \in \text{safe}(R, \epsilon, L, \lambda, \omega)$.*

The proof of Proposition 4.1 can be found in Appendix D.4 (see Theorem D.21) and is again an application of Berge’s theorem (Berge, 1963). Note that the regret bound L is defined as the regret of the reference policy. This makes intuitive sense, as regularized policy optimization constrains the policy under optimization $\hat{\pi}$ to not deviate too strongly from the reference policy π_{ref} , which will also constrain the regret of $\hat{\pi}$ to stay close to the regret of π_{ref} . Under the conditions of Proposition 4.1, the regret of π_{ref} serves as an upper regret bound because for small enough ϵ the learned reward \hat{R} and the true reward R are close enough such that maximizing \hat{R} also improve reward with respect to R . Furthermore, we note that it is also possible to derive a version of the theorem in which the expected error in rewards is replaced by a KL divergence in choice probabilities, similar to Proposition D.14, by combining the arguments in that proposition with the arguments in Berge’s theorem. A full formulation and proof of the result can be found in Theorem D.22.

Similar to Proposition 3.1, Proposition 4.1 does not guarantee the existence of a universal ϵ such that all data distributions D are in $\text{safe}(R, \epsilon, L, \lambda, \omega)$. In our next result, we show that such an ϵ does not exist, since for each ϵ , there is a nontrivial set of data distributions that allows for error-regret mismatch to occur:

Theorem 4.2. *Let $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an arbitrary MDP, $\lambda \in (0, \infty)$, $L \in (0, 1)$, and $\omega : \Pi \rightarrow \mathbb{R}$ be a regularization function. Furthermore, let π_* be a deterministic worst-case policy for R , meaning that $\text{Reg}^R(\pi_*) = 1$. Let $C := C(\mathcal{M}, \pi_*, L, \lambda, \omega) < \infty$ be the constant defined in Equation (98) in the appendix. Let $\epsilon > 0$. Then for all data distributions $D \in \Delta(\mathcal{S} \times \mathcal{A})$ with*

$$D(\text{supp } D^{\pi_*}) \leq \frac{\epsilon}{1 + C}, \quad (4)$$

we have $D \in \text{unsafe}(R, \epsilon, L, \lambda, \omega)$.

The proof of Theorem 4.2 can be found in Appendix C.5 (see Theorem C.38). The general idea is as follows: To prove that D is unsafe, define \hat{R} to be equal to R outside of $\text{supp } D^{\pi_*}$, and very large in $\text{supp } D^{\pi_*}$. If it is sufficiently large in this region, then regularized optimization leads to a policy $\hat{\pi}$

¹E.g., if $\pi_{\text{ref}}(a | s) > 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $\omega(\pi) := \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}})$, then the minimum is π_{ref} .

with $\text{Reg}^R(\hat{\pi}) \geq L$. Finally, the condition that $D(\text{supp } D^{\pi^*}) \leq \frac{\epsilon}{1+C}$ ensures that \hat{R} has a reward error bounded by ϵ .

Note that Theorem 4.2 is very general and covers a large class of different regularization methods. In Corollary C.40 we provide a specialized result for the case of KL -regularized policy optimization, and in Section 6 we investigate error-regret mismatch in the RLHF framework.

5 GENERALIZATION OF THE ERROR MEASUREMENT

Our results have so far expressed the error of the learned reward \hat{R} in terms of Equation (1), i.e., in terms of the expected error of individual transitions. In this section, we show that many common reward learning training objectives can be upper-bounded in terms of the expected error metric defined in Equation (1). This in turn means that our negative results generalize to reward learning algorithms that use these other training objectives. We state all upper bounds for MDPs with finite time horizon T (but note that these results directly generalize to MDPs with infinite time horizon by taking the limit of $T \rightarrow \infty$).

In the finite horizon setting, trajectories are defined as a finite list of states and actions: $\xi = s_0, a_0, s_1, \dots, a_{T-1}$. We use Ξ for the set of all trajectories of length T . As in the previous sections, $G: \Xi \rightarrow \mathbb{R}$ denotes the trajectory return function, defined as $G(\xi) = \sum_{t=0}^{T-1} \gamma^t \cdot R(s_t, a_t)$. We start by showing that low expected error in transitions implies low expected error in trajectory returns:

Proposition 5.1. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and its resulting data distribution $D^\pi = \frac{1-\gamma}{1-\gamma^T} \cdot \eta^\pi$ and a second reward function $\hat{R}: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected difference in trajectory evaluation as follows:*

$$\mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|] \leq \frac{1-\gamma^T}{1-\gamma} \cdot \mathbb{E}_{(s,a) \sim D^\pi} [|R(s,a) - \hat{R}(s,a)|].$$

The proof of Proposition 5.1 can be found in Appendix C.4.1 (see Proposition C.24). Furthermore, a low expected error of trajectory returns implies a low expected error of choice distributions (a distance metric commonly used as the loss in RLHF (Christiano et al., 2017)). Namely, given a reward function R , define the probability of trajectory ξ_1 being preferred over ξ_2 to be $p_R(\xi_1 \succ \xi_2) = \sigma(G_R(\xi_1) - G_R(\xi_2)) = \frac{\exp(G_R(\xi_1))}{\exp(G_R(\xi_1)) + \exp(G_R(\xi_2))}$. We then have:

Proposition 5.2. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and a second reward function $\hat{R}: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected KL divergence over trajectory preference distributions as follows:*

$$\mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [\mathbb{D}_{KL}(p_R(\cdot | \xi_1, \xi_2) || p_{\hat{R}}(\cdot | \xi_1, \xi_2))] \leq 2 \cdot \mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|].$$

The proof of Proposition 5.2 can be found in Appendix C.4.1 (see Proposition C.25).

Finally, in some RLHF scenarios, for example in RLHF with prompt-response pairs, one prefers to only compare trajectories with a common starting state. In the following proposition, we upper-bound the expected error of choice distributions with trajectories that share a common starting state by the expected error of choice distributions with arbitrary trajectories:

Proposition 5.3. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and a second reward function $\hat{R}: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected KL divergence of preference distributions over trajectories with a common starting state as follows:*

$$\mathbb{E}_{\substack{s_0 \sim \mu_0, \\ \xi_1, \xi_2 \sim \pi(s_0)}} [\mathbb{D}_{KL}(p_R(\cdot | \xi_1, \xi_2) || p_{\hat{R}}(\cdot | \xi_1, \xi_2))] \leq \frac{\mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [\mathbb{D}_{KL}(p_R(\cdot | \xi_1, \xi_2) || p_{\hat{R}}(\cdot | \xi_1, \xi_2))]}{\min_{s' \in \mathcal{S}, \mu_0(s') > 0} \mu_0(s')}.$$

The proof of Proposition 5.3 can be found in Appendix C.4.1 (see Proposition C.26).

6 ERROR-REGRET MISMATCH IN RLHF

In this section we use the generalization results from Section 5 to extend our results to reinforcement learning from human feedback (RLHF). We provide more general results about the class of KL -regularized optimization policy optimization methods in Appendix C.4.5.

486 RLHF, especially in the context of large language models, is usually modeled in a *contextual bandit*
 487 setting (Ziegler et al., 2019; Stiennon et al., 2020; Bai et al., 2022; Ouyang et al., 2022; Rafailov
 488 et al., 2023). A *contextual bandit* $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ is defined by a set of states \mathcal{S} , a set of actions \mathcal{A} , a
 489 data distribution $\mu_0 \in \Delta(\mathcal{S})$, and a reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The goal is to learn a policy
 490 $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that maximizes the expected return $J(\pi) = \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [R(s, a)]$. In the context
 491 of language models, \mathcal{S} is usually called the set of *prompts* or *contexts*, and \mathcal{A} the set of *responses*.

492 We state the following theorem using a more precise version of Definition 2.1 tailored to the
 493 RLHF setting. In particular, we replace the similarity metric (property 1. of Definition 2.1)
 494 with the expected similarity in choice probabilities. A precise mathematical definition can be
 495 found in Appendix C.4.3. We denote the resulting sets of safe- and unsafe data distributions by
 496 $\text{safe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}}))$ and $\text{unsafe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}}))$.

497 By making use of the specifics of this setting, we can derive more interpretable and stronger results.
 498 In particular, we specify a set of reference distributions for which performing KL-regularized policy
 499 optimization allows for error-regret mismatch to occur.

500 **Theorem 6.1.** *Let $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ be a contextual bandit. Given $L \in [0, 1)$, we define for every state
 501 $s \in \mathcal{S}$ the reward threshold: $R_L(s) := (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a)$. Lastly,
 502 let $\pi_{\text{ref}} : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ be an arbitrary reference policy for which it holds that for every $(s, a) \in \mathcal{S} \times \mathcal{A}$,
 503 $\pi_{\text{ref}}(a|s) > 0$, and there exists at least one action $a_s \in \mathcal{A}$ such that $R(s, a_s) < R_L(s)$ and $\pi_{\text{ref}}(a_s|s)$
 504 satisfies the following inequality:*

$$505 \pi_{\text{ref}}(a_s|s) \leq \frac{(R_L(s) - R(s, a_s))}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{4 \cdot \lambda^2}.$$

506 Let $D^{\text{ref}}(s, a) := \mu_0(s) \cdot \pi_{\text{ref}}(a|s)$. Then $D^{\text{ref}} \in \text{unsafe}^{\text{RLHF}}(R, 2 \cdot \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}}))$

507 The proof of Theorem 6.1 can be found in Appendix C.4.4 (see Propositions C.31 and C.32). We
 508 expect the conditions on the reference policy π_{ref} to be likely to hold in real-world cases as the
 509 number of potential actions (or responses) is usually very large, and language models typically assign
 510 a large portion of their probability mass to only a tiny fraction of all responses. This means that for
 511 every state/prompt s , a huge majority of actions/responses a have a very small probability $\pi_{\text{ref}}(a|s)$.

512 7 DISCUSSION

513 We have contributed to building up the foundations for the learning theory of general reward learning
 514 in arbitrary MDPs by studying the relationship between the expected error of a learned reward function
 515 on some data distribution and the extent to which optimizing that reward function is guaranteed
 516 to produce a policy with low regret according to the true reward function. We showed that as the
 517 expected error ϵ of a reward model \hat{R} goes to zero, the worst-case regret of a policy that is optimal
 518 under \hat{R} (with or without regularization) also goes to zero (Propositions 3.1 and 4.1). However, in
 519 Proposition 3.2 we also showed that ϵ , in general, must be extremely small to ensure that \hat{R} 's optimal
 520 policies have a low worst-case regret. In particular, this value depends on the smallest probability that
 521 the data distribution D assigns to any transition in the underlying MDP, which means that it shrinks
 522 very quickly for large MDPs. Consequently, there exists no ϵ that can universally ensure low regret.

523 More generally, low expected error does *not* ensure low regret for all realistic data distributions
 524 (Proposition 3.3, Theorem 4.2 and Theorem 6.1). We refer to this phenomenon as *error-regret*
 525 *mismatch*. This is due to policy optimization (typically) involving a *distributional shift* from the
 526 data distribution that is used to train the reward model; a reward model that is accurate on the
 527 data distribution may fail to be accurate after this distributional shift. Moreover, we find evidence
 528 that some MDPs with very large action spaces do not allow for *any* safe data distributions relative
 529 to a reasonable reward model error and desired regret bound (Corollary 3.4). We also showed
 530 that our results generalize to various different data sources, such as preferences over trajectories
 531 (Propositions 5.2 and 5.3) and trajectory scores (Proposition 5.1), supporting the conclusion that this
 532 issue is a fundamental problem of reward learning.

533 Lastly, for the case of unregularized optimization, we derive a set of *necessary and sufficient*
 534 conditions that allow us to determine the set of safe and unsafe data distributions for arbitrary MDPs,
 535 thereby completely answering the question of when exactly a data distribution is safe (Theorem 3.5).

Our results highlight the challenge of deriving useful PAC-like generalization bounds for current reward learning algorithms. While there do exist bounds (Nika et al., 2024; Cen et al., 2024), they depend on some form of data coverage of (bad) policies. As we have shown, in practical situations, we should expect the coverage to be so low that the regret will be large. **Our results highlight the challenge of deriving useful PAC-like generalization bounds for current reward learning algorithms. While this is possible (and has been done, see (Nika et al., 2024; Cen et al., 2024)), we showed that realistic bounds on the error in the learned reward function do not provide meaningful guarantees. By focusing on the propagation of reward function error to regret in policy optimization, our work provides an insightful analysis that disentangles a key obstacle specific to reward learning from classical learning theory challenges.**

Our results also highlight the importance of researching ways for evaluating reward functions using methods other than evaluating them on a test set, e.g. by using interpretability methods (Michaud et al., 2020; Jenner & Gleave, 2022) or finding better ways to quantify reward function distance (Gleave et al., 2020; Skalse et al., 2024)). **These are largely unsolved research efforts that would profit from further engagement.**

7.1 LIMITATIONS AND FUTURE WORK

Our work focuses on the question of whether there *exists* a reward model \hat{R} that is compatible with the true reward function on a data distribution, such that there *exists* a policy $\hat{\pi}$ that is optimal under \hat{R} , but which has high regret. In practice, it may be that the inductive bias of the reward learning algorithm or the policy optimization algorithm avoids the most pathological cases. Our analysis could therefore be extended by attempting to take the inductive bias into account. Furthermore, our analyses assume that we are able to find optimal policies, but in practice, this is rarely the case. Generalizing our results to non-optimal policies therefore constitutes an important direction for further research. **Finally, one could attempt to analyze the *likelihood* of a high-regret training outcome of reward learning and policy optimization instead of analyzing the worst-case.**

Furthermore, it is important to theoretically analyze improved reward learning and policy optimization procedures. There is already some empirical work on using reward model ensembles (Coste et al., 2023) or weight averaged reward models (Ramé et al., 2024) to overcome problems of reward model overoptimization. In the special case of multi-armed bandits, iterated data-smoothing has been proposed and analyzed theoretically and empirically (Zhu et al., 2024). Very recent work also considers learning reward models on online data for mitigating distribution shifts and thus reward overoptimization (Lang et al., 2024a) or even theoretically analyzes such a setting for the special case of linear reward functions (Song et al., 2024). We hope that a careful theoretical analysis of all these settings in similar generality as our work can identify reliable ways to improve upon the “theoretical baseline” established by our work.

In addition to improving the theory and practice of reward learning itself, there are other ways to improve the safety of the resulting policies after training. We are excited about efforts to evaluate policies for dangerous capabilities (Phuong et al., 2024), red-teaming (Perez et al., 2022), safety cases (Clymer et al., 2024), shields (Alshiekh et al., 2018), and a numerous suite of other approaches (Anwar et al., 2024).

Moreover, there are numerous opportunities to identify more necessary and/or sufficient conditions when a data distribution (dis)allows error-regret mismatch. In general, it would be interesting to find more interpretable and practical conditions that guarantee a data distribution is safe or unsafe, i.e., conditions that do not rely on knowledge about the true reward function or the transition distribution.

For the purposes of communicating our paper updates in the rebuttal, we made additions to the paper that address some of the reviewer’s concerns. They temporarily increase the page number above the limit. We will make sure to fit everything within the page limit for the camera-ready version.

REFERENCES

- 594
595
596 Alekh Agarwal, Miroslav Dudík, Satyen Kale, John Langford, and Robert Schapire. Contextual
597 bandit learning with predictable rewards. In *Artificial Intelligence and Statistics*, pp. 19–26. PMLR,
598 2012.
- 599 Anurag Ajay, Abhishek Gupta, Dibya Ghosh, Sergey Levine, and Pulkit Agrawal. Distributionally
600 adaptive meta reinforcement learning. *Advances in Neural Information Processing Systems*, 35:
601 25856–25869, 2022.
- 602 Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and
603 Ufuk Topcu. Safe Reinforcement Learning via Shielding. *Proceedings of the AAAI Conference
604 on Artificial Intelligence*, 32(1), Apr. 2018. doi: 10.1609/aaai.v32i1.11797. URL <https://ojs.aaai.org/index.php/AAAI/article/view/11797>.
- 605 Anthropic. Introducing Claude. [https://www.anthropic.com/index/
606 introducing-claude](https://www.anthropic.com/index/introducing-claude), 2023. Accessed: 2023-09-05.
- 607 Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase,
608 Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, Benjamin L. Edelman,
609 Zhaowei Zhang, Mario Günther, Anton Korinek, Jose Hernandez-Orallo, Lewis Hammond, Eric
610 Bigelow, Alexander Pan, Lauro Langosco, Tomasz Korbak, Heidi Zhang, Ruiqi Zhong, Seán Ó
611 hÉigeartaigh, Gabriel Recchia, Giulio Corsi, Alan Chan, Markus Anderljung, Lilian Edwards, Alek-
612 sandar Petrov, Christian Schroeder de Witt, Sumeet Ramesh Motwan, Yoshua Bengio, Danqi Chen,
613 Philip H. S. Torr, Samuel Albanie, Tegan Maharaj, Jakob Foerster, Florian Tramer, He He, Atoosa
614 Kasirzadeh, Yejin Choi, and David Krueger. Foundational Challenges in Assuring Alignment and
615 Safety of Large Language Models, 2024. URL <https://arxiv.org/abs/2404.09932>.
- 616 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain,
617 Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with
618 reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- 619 Andrea Bajcsy, Dylan P Losey, Marcia K O’malley, and Anca D Dragan. Learning robot objectives
620 from physical human interaction. In *Conference on robot learning*, pp. 217–226. PMLR, 2017.
- 621 Claude Berge. *Topological Spaces: Including a Treatment of Multi-valued Functions, Vector
622 Spaces and Convexity*. Macmillan, 1963. URL [https://books.google.nl/books?id=
623 0QJRAAAAMAAJ](https://books.google.nl/books?id=0QJRAAAAMAAJ).
- 624 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. The method
625 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 626 Daniel S Brown and Scott Niekum. Deep Bayesian reward learning from preferences. *arXiv preprint
627 arXiv:1912.04472*, 2019.
- 628 Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier
629 Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems
630 and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint
631 arXiv:2307.15217*, 2023.
- 632 Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale
633 Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified
634 approach to online and offline rlhf. *arXiv preprint arXiv:2405.19320*, 2024.
- 635 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
636 reinforcement learning from human preferences. *Advances in neural information processing
637 systems*, 30, 2017.
- 638 Joshua Clymer, Nick Gabrieli, David Krueger, and Thomas Larsen. Safety Cases: How to Justify the
639 Safety of Advanced AI Systems, 2024. URL <https://arxiv.org/abs/2403.10462>.
- 640 Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization
641 in reinforcement learning. In *International conference on machine learning*, pp. 1282–1289. PMLR,
642 2019.

- 648 Thomas Coste, Usman Anwar, Robert Kirk, and David Krueger. Reward model ensembles help
649 mitigate overoptimization. *arXiv preprint arXiv:2310.02743*, 2023.
- 650
- 651 Jesse Farebrother, Marlos C Machado, and Michael Bowling. Generalization and regularization in
652 dqn. *arXiv preprint arXiv:1810.00123*, 2018.
- 653
- 654 Dylan J Foster, Alexander Rakhlin, David Simchi-Levi, and Yunzong Xu. Instance-dependent
655 complexity of contextual bandits and reinforcement learning: A disagreement-based perspective.
656 *arXiv preprint arXiv:2010.03104*, 2020.
- 657
- 658 Ted Fujimoto, Joshua Suetterlein, Samrat Chatterjee, and Auroop Ganguly. Assessing the impact of
659 distribution shift on reinforcement learning performance. *arXiv preprint arXiv:2402.03590*, 2024.
- 660
- 661 Leo Gao, John Schulman, and Jacob Hilton. Scaling laws for reward model overoptimization. In
662 *International Conference on Machine Learning*, pp. 10835–10866. PMLR, 2023.
- 663
- 664 Google Gemini Team. Gemini: A Family of Highly Capable Multimodal Mod-
665 els. https://storage.googleapis.com/deepmind-media/gemini/gemini_1_report.pdf, 2023. Accessed: 2023-12-11.
- 666
- 667 Adam Gleave, Michael Dennis, Shane Legg, Stuart Russell, and Jan Leike. Quantifying differences
668 in reward functions. *arXiv preprint arXiv:2006.13900*, 2020.
- 669
- 670 Charles AE Goodhart. *Problems of monetary management: the UK experience*. Springer, 1984.
- 671
- 672 Christopher A Hennessy and Charles AE Goodhart. Goodhart’s law and machine learning: a structural
673 perspective. *International Economic Review*, 64(3):1075–1086, 2023.
- 674
- 675 Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward
676 learning from human preferences and demonstrations in Atari. In *Proceedings of the 32nd
677 International Conference on Neural Information Processing Systems*, volume 31, pp. 8022–8034,
678 Montréal, Canada, 2018. Curran Associates, Inc., Red Hook, NY, USA.
- 679
- 680 Erik Jenner and Adam Gleave. Preprocessing reward functions for interpretability, 2022.
- 681
- 682 Hong Jun Jeon, Smitha Milli, and Anca Dragan. Reward-rational (implicit) choice: A unifying
683 formalism for reward learning. *Advances in Neural Information Processing Systems*, 33:4415–4426,
684 2020.
- 685
- 686 Xiang Ji, Huazheng Wang, Minshuo Chen, Tuo Zhao, and Mengdi Wang. Provable benefits of policy
687 learning from human preferences in contextual bandit problems. *arXiv preprint arXiv:2307.12975*,
688 2023.
- 689
- 690 Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In
691 *International Conference on Machine Learning*, pp. 5084–5096. PMLR, 2021.
- 692
- 693 Jacek Karwowski, Oliver Hayman, Xingjian Bai, Klaus Kiendlhofer, Charlie Griffin, and Joar Skalse.
694 Goodhart’s Law in Reinforcement Learning. *arXiv preprint arXiv:2310.09144*, 2023.
- 695
- 696 Michael J. Kearns and Umesh Vazirani. *An Introduction to Computational Learning Theory*. The
697 MIT Press, 08 1994. ISBN 9780262276863. doi: 10.7551/mitpress/3897.001.0001. URL
698 <https://doi.org/10.7551/mitpress/3897.001.0001>.
- 699
- 700 Kihyun Kim, Jiawei Zhang, Pablo A Parrilo, and Asuman Ozdaglar. A unified linear programming
701 framework for offline reward learning from human demonstrations and feedback. *arXiv preprint
arXiv:2405.12421*, 2024.
- 702
- 703 Victoria Krakovna. Specification gaming: The flip side of Ai Ingenu-
704 ity, Apr 2020. URL [https://deepmind.google/discover/blog/
705 specification-gaming-the-flip-side-of-ai-ingenuity/](https://deepmind.google/discover/blog/specification-gaming-the-flip-side-of-ai-ingenuity/).
- 706
- 707 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
708 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.

- 702 Thomas Kwa, Drake Thomas, and Adrià Garriga-Alonso. Catastrophic goodhart: regularizing
703 rlhf with kl divergence does not mitigate heavy-tailed reward misspecification. *arXiv preprint*
704 *arXiv:2407.14503*, 2024.
- 705 Cassidy Laidlaw, Shivam Singhal, and Anca Dragan. Preventing reward hacking with occupancy
706 measure regularization. *arXiv preprint arXiv:2403.03185*, 2024.
- 707 Hao Lang, Fei Huang, and Yongbin Li. Fine-Tuning Language Models with Reward Learning on
708 Policy. *arXiv preprint arXiv:2403.19279*, 2024a.
- 709 Leon Lang, Davis Foote, Stuart Russell, Anca Dragan, Erik Jenner, and Scott Emmons. When Your
710 AIs Deceive You: Challenges with Partial Observability of Human Evaluators in Reward Learning.
711 *arXiv preprint arXiv:2402.17747*, 2024b.
- 712 Haoyang Li, Xin Wang, Ziwei Zhang, and Wenwu Zhu. Out-of-distribution generalization on graphs:
713 A survey. *arXiv preprint arXiv:2202.07987*, 2022.
- 714 Ying Li, Xingwei Wang, Rongfei Zeng, Praveen Kumar Donta, Ilir Murturi, Min Huang, and
715 Schahram Dustdar. Federated domain generalization: A survey. *arXiv preprint arXiv:2306.01334*,
716 2023.
- 717 Jiashuo Liu, Zheyang Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards
718 out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.
- 719 Miao Lu, Han Zhong, Tong Zhang, and Jose Blanchet. Distributionally robust reinforcement learning
720 with interactive data collection: Fundamental hardness and near-optimal algorithm. *arXiv preprint*
721 *arXiv:2404.03578*, 2024.
- 722 Khushdeep Singh Mann, Steffen Schneider, Alberto Chiappa, Jin Hwa Lee, Matthias Bethge, Alexan-
723 der Mathis, and Mackenzie W Mathis. Out-of-distribution generalization of internal models is
724 correlated with reward. In *Self-Supervision for Reinforcement Learning Workshop-ICLR*, volume
725 2021, 2021.
- 726 Viraj Mehta, Vikramjeet Das, Ojash Neopane, Yijia Dai, Ilija Bogunovic, Jeff Schneider, and Willie
727 Neiswanger. Sample efficient reinforcement learning from human feedback via active exploration.
728 *OpenReview*, 2023.
- 729 Eric J. Michaud, Adam Gleave, and Stuart Russell. Understanding learned reward functions, 2020.
- 730 Andrew Y Ng and Stuart Russell. Algorithms for inverse reinforcement learning. In *Proceedings of*
731 *the Seventeenth International Conference on Machine Learning*, volume 1, pp. 663–670, Stanford,
732 California, USA, 2000. Morgan Kaufmann Publishers Inc.
- 733 Andrew Y Ng, Stuart Russell, et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1,
734 pp. 2, 2000.
- 735 Andi Nika, Debmalya Mandal, Parameswaran Kamalaruban, Georgios Tzannetos, Goran Radanović,
736 and Adish Singla. Reward Model Learning vs. Direct Policy Optimization: A Comparative
737 Analysis of Learning from Human Preferences. *arXiv preprint arXiv:2403.01857*, 2024.
- 738 OpenAI. Introducing ChatGPT. <https://openai.com/blog/chatgpt>, 2022. Accessed:
739 2024-02-06.
- 740 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
741 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
742 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
743 27744, 2022.
- 744 Malayandi Palan, Nicholas Charles Landolfi, Gleb Shevchuk, and Dorsa Sadigh. Learning reward
745 functions by integrating human demonstrations and preferences. In *Proceedings of Robotics:*
746 *Science and Systems*, Freiburg im Breisgau, Germany, June 2019. doi: 10.15607/RSS.2019.XV.023.
- 747 Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese,
748 Nat McAleese, and Geoffrey Irving. Red Teaming Language Models with Language Models, 2022.
749 URL <https://arxiv.org/abs/2202.03286>.

- 756 Mary Phuong, Matthew Aitchison, Elliot Catt, Sarah Cogan, Alexandre Kaskasoli, Victoria Krakovna,
757 David Lindner, Matthew Rahtz, Yannis Assael, Sarah Hodkinson, Heidi Howard, Tom Lieberum,
758 Ramana Kumar, Maria Abi Raad, Albert Webson, Lewis Ho, Sharon Lin, Sebastian Farquhar, Mar-
759 cus Hutter, Gregoire Deletang, Anian Ruoss, Seliem El-Sayed, Sasha Brown, Anca Dragan, Rohin
760 Shah, Allan Dafoe, and Toby Shevlane. Evaluating Frontier Models for Dangerous Capabilities,
761 2024. URL <https://arxiv.org/abs/2403.13793>.
- 762 Martin L Puterman. Markov Decision Processes: Discrete Stochastic Dynamic Programming, 1994.
763
- 764 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea
765 Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv*
766 *preprint arXiv:2305.18290*, 2023.
- 767 Alexandre Ramé, Nino Vieillard, Léonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier
768 Bachem, and Johan Ferret. Warm: On the benefits of weight averaged reward models. *arXiv*
769 *preprint arXiv:2401.12187*, 2024.
- 770 R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science &
771 Business Media, 2009.
- 772
- 773 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
774 prediction to no-regret online learning. In *Proceedings of the fourteenth international conference*
775 *on artificial intelligence and statistics*, pp. 627–635. JMLR Workshop and Conference Proceedings,
776 2011.
- 777 Andreas Schleginhausen and Maryam Kamgarpour. Identifiability and generalizability in constrained
778 inverse reinforcement learning. In *International Conference on Machine Learning*, pages=30224–
779 30251. PMLR, 2023.
- 780
- 781 Rohin Shah, Dmitrii Krasheninnikov, Jordan Alexander, Pieter Abbeel, and Anca Dragan. Preferences
782 Implicit in the State of the World. *arXiv e-prints*, art. arXiv:1902.04198, February 2019. doi:
783 10.48550/arXiv.1902.04198.
- 784 Joar Skalse and Alessandro Abate. Misspecification in inverse reinforcement learning, 2023.
785
- 786 Joar Skalse and Alessandro Abate. Quantifying the sensitivity of inverse reinforcement learning to
787 misspecification, 2024.
- 788 Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing
789 reward gaming. *Advances in Neural Information Processing Systems*, 35:9460–9471, 2022.
- 790
- 791 Joar Skalse, Lucy Farnik, Sumeet Ramesh Motwani, Erik Jenner, Adam Gleave, and Alessandro
792 Abate. Starc: A general framework for quantifying differences between reward functions, 2024.
- 793 Joar Max Viktor Skalse, Matthew Farrugia-Roberts, Stuart Russell, Alessandro Abate, and Adam
794 Gleave. Invariance in policy optimisation and partial identifiability in reward learning. In *Internat-*
795 *ional Conference on Machine Learning*, pp. 32033–32058. PMLR, 2023.
- 796
- 797 Yuda Song, Gokul Swamy, Aarti Singh, J. Andrew Bagnell, and Wen Sun. The Importance of Online
798 Data: Understanding Preference Fine-tuning via Coverage, 2024. URL <https://arxiv.org/abs/2406.01462>.
799
- 800 Richard Stanley. Chapter 1: Basic Definitions, the Intersection Poset and the Characteristic
801 Polynomial. In *Combinatorial Theory: Hyperplane Arrangements—MIT Course No. 18.315*.
802 MIT OpenCourseWare, Cambridge MA, 2024. URL [https://ocw.mit.edu/courses/](https://ocw.mit.edu/courses/18-315-combinatorial-theory-hyperplane-arrangements-fall-2004/pages/lecture-notes/)
803 [18-315-combinatorial-theory-hyperplane-arrangements-fall-2004/](https://ocw.mit.edu/courses/18-315-combinatorial-theory-hyperplane-arrangements-fall-2004/pages/lecture-notes/)
804 [pages/lecture-notes/](https://ocw.mit.edu/courses/18-315-combinatorial-theory-hyperplane-arrangements-fall-2004/pages/lecture-notes/). MIT OpenCourseWare.
- 805 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
806 Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in*
807 *Neural Information Processing Systems*, 33:3008–3021, 2020.
- 808
- 809 Marilyn Strathern. ‘Improving ratings’: audit in the British University system. *European review*, 5
(3):305–321, 1997.

- 810 Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, second
811 edition, 2018. ISBN 9780262352703.
- 812
- 813 Jeremy Tien, Jerry Zhi-Yang He, Zackory Erickson, Anca D Dragan, and Daniel S Brown. Causal
814 confusion and reward misidentification in preference-based reward learning. *arXiv preprint*
815 *arXiv:2204.06601*, 2022.
- 816 Hsiao-Yu Tung, Adam W Harley, Liang-Kang Huang, and Katerina Fragkiadaki. Reward learning
817 from narrated demonstrations. In *Proceedings: 2018 IEEE/CVF Conference on Computer Vision*
818 *and Pattern Recognition (CVPR)*, pp. 7004–7013, Salt Lake City, Utah, USA, June 2018. IEEE
819 Computer Society, Los Alamitos, CA, USA. doi: 10.1109/CVPR.2018.00732.
- 820 Robert J Vanderbei. Linear programming: foundations and extensions. *Journal of the Operational*
821 *Research Society*, 49(1):94–94, 1998.
- 822
- 823 Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun
824 Zeng, and S Yu Philip. Generalizing to unseen domains: A survey on domain generalization. *IEEE*
825 *transactions on knowledge and data engineering*, 35(8):8052–8072, 2022.
- 826 Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre Alvisé-Rebuffi, Ira Ktena, Krishnamurthy
827 Dvijotham, and Taylan Cemgil. A fine-grained analysis on distribution shift. *arXiv preprint*
828 *arXiv:2110.11328*, 2021.
- 829
- 830 Christian Wirth, Riad Akrou, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-
831 based reinforcement learning methods. *Journal of Machine Learning Research*, 18(136):1–46,
832 2017.
- 833 Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang.
834 Iterative preference learning from human feedback: Bridging theory and practice for rlhf under
835 kl-constraint. In *Forty-first International Conference on Machine Learning*, 2024.
- 836 Jee Seok Yoon, Kwansoek Oh, Yooseung Shin, Maciej A Mazurowski, and Heung-Il Suk. Domain
837 Generalization for Medical Image Analysis: A Survey. *arXiv preprint arXiv:2310.08598*, 2023.
- 838
- 839 Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable Offline
840 Preference-Based Reinforcement Learning. In *The Twelfth International Conference on Learning*
841 *Representations*, 2023.
- 842 Amy Zhang, Nicolas Ballas, and Joelle Pineau. A dissection of overfitting and generalization in
843 continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*, 2018.
- 844
- 845 Xuezhou Zhang, Yiding Chen, Xiaojin Zhu, and Wen Sun. Corruption-robust offline reinforcement
846 learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 5757–5773.
847 PMLR, 2022.
- 848 Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A
849 survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4396–4415, 2022.
- 850
- 851 Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human
852 feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*,
853 pp. 43037–43067. PMLR, 2023.
- 854
- 855 Banghua Zhu, Michael I Jordan, and Jiantao Jiao. Iterative data smoothing: Mitigating reward
856 overfitting and overoptimization in rlhf. *arXiv preprint arXiv:2401.16335*, 2024.
- 857
- 858 Simon Zhuang and Dylan Hadfield-Menell. Consequences of misaligned AI. In *Proceedings of the*
859 *34th International Conference on Neural Information Processing Systems, NIPS’20*, pp. 15763–
860 15773, Red Hook, NY, USA, December 2020. Curran Associates Inc. ISBN 978-1-71382-954-6.
- 861
- 862 Simon Zhuang and Dylan Hadfield-Menell. Consequences of misaligned AI. *Advances in Neural*
863 *Information Processing Systems*, 33:15763–15773, 2020.
- 864
- 865 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul
866 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv*
867 *preprint arXiv:1909.08593*, 2019.

APPENDIX

This appendix develops the theory outlined in the main paper in a self-contained and complete way, including all proofs. In Appendix B, we present the setup of all concepts and the problem formulation, as was already contained in the main paper. In Appendix C, we present all “negative results”. Conditional on an error threshold in the reward model, these results present conditions for the data distribution that allow reward models to be learned that allow for error-regret mismatch. That section also contains Theorem C.16 which is an equivalent condition for the absence of error-regret mismatch but could be considered a statement about error-regret mismatch by negation. In Appendix D, we present sufficient conditions for *safe optimization* in several settings. Typically, this boils down to showing that given a data distribution, a *sufficiently small* error in the reward model guarantees that its optimal policies have low regret.

CONTENTS OF THE APPENDIX

A	Extended related work	18
B	Introduction	19
	B.1 Preliminaries	19
	B.2 Problem formalization	19
C	Existence of error-regret mismatch	20
	C.1 Assumptions	20
	C.2 Intuitive unregularized existence statement	20
	C.3 General existence statements	22
	C.3.1 More interpretable statement	26
	C.3.2 Deriving the conditions on D	28
	C.3.3 Algorithm to compute the conditions on D	33
	C.3.4 Working example of computing matrix M	34
	C.3.5 Building up on Theorem 3.5	35
	C.4 Existence of negative results in the RLHF setting	35
	C.4.1 Generalization of the error measurement	35
	C.4.2 RLHF bandit formulation	38
	C.4.3 Safe and unsafe data distributions for RLHF	39
	C.4.4 Negative results	40
	C.4.5 Another negative result with regularization	47
	C.5 A regularized negative result for general MDPs	49
D	Requirements for safe optimization	56
	D.1 Applying Berge’s maximum theorem	56
	D.2 Elementary proof of a regret bound	59
	D.3 Safe optimization via approximated choice probabilities	62
	D.4 Positive result for regularized RLHF	65

A EXTENDED RELATED WORK

Reward Learning Reward learning is a key concept in reinforcement learning that involves learning the reward function for complex tasks with latent and difficult-to-specify reward functions. Many methods have been developed to incorporate various types of human feedback into the reward learning process (Wirth et al., 2017; Ng et al., 2000; Bajcsy et al., 2017; Jeon et al., 2020).

Challenges in Reward Learning Reward learning presents several challenges (Casper et al., 2023; Lang et al., 2024b; Skalse & Abate, 2023; 2024), such as *reward misgeneralization*, where the reward model learns a different reward function that performs well on in-distribution data but differs strongly on out-of-distribution data (Skalse et al., 2023). This can lead to unintended consequences in real-world applications.

Reward misgeneralization can also result in *reward hacking* (Krakovna, 2020), a consequence of Goodhart’s law (Goodhart, 1984; Zhuang & Hadfield-Menell, 2020; Hennessy & Goodhart, 2023; Strathern, 1997; Karwowski et al., 2023). Reward hacking has been extensively studied both theoretically (Skalse et al., 2022; 2024; Zhuang & Hadfield-Menell, 2020) and empirically (Zhang et al., 2018; Farebrother et al., 2018; Cobbe et al., 2019; Krakovna, 2020; Gao et al., 2023; Tien et al., 2022).

How we complement prior work Over the past few years, several works have observed and investigated distribution shifts in RL empirically (Wiles et al., 2021; Fujimoto et al., 2024; Mann et al., 2021; Ajay et al., 2022; Kumar et al., 2020). These studies provide valuable insights into the impact of distribution shifts on RL performance, but there remain gaps in our understanding of this phenomenon.

We are interested in theoretically analyzing the impact that the initial *data distribution*, used to train the reward model, has on the regret of the final policy. In particular, we assume arbitrary MDPs and investigate the regret a policy might attain for a worst-case reward model, that might be trained during reward learning. Our work thereby nicely extends and complements theoretical work that investigates policy regret under a worst-case MDP in settings such as imitation learning (Ross et al., 2011), offline RL (Kim et al., 2024; Jin et al., 2021; Zhang et al., 2022), and other RL settings (Lu et al., 2024; Laidlaw et al., 2024; Kwa et al., 2024)

Furthermore, we try to remain as general in our results as possible by letting all our results in Sections 3 and 4 hold for arbitrary MDPs without any additional restrictions. In contrast, (Jin et al., 2021; Zhang et al., 2022; Nika et al., 2024) choose to investigate results in linear MDPs, whereas (Zhu et al., 2024) focus their analysis to multi-armed bandits.

In terms of new results, we develop precise conditions (i.e., necessary and sufficient) for when exactly a given data distribution is safe for some fixed MDP (Theorem 3.5) for unregularized policy optimization. To the best of our knowledge, we are the first to develop such results and point out the surprisingly regular shape of the set of safe data distributions. We hope this can act as a “theoretical baseline” for theoretical investigations of improved reward learning methods and their “safe starting conditions”.

Lastly, we demonstrate that without additional assumptions, a wide range of reward learning methods are not guaranteed to be safe, i.e., there exist reasonable conditions under which RL with reward learning might yield a policy that behaves very badly. Importantly, straightforward fixes such as policy regularization do not fix this issue (Theorem 4.2) and we show in Section 6 that our results directly apply to the standard RLHF setting as well. This means that the most widely used LLM safety technique is not safe, without additional assumptions. In contrast to the works Zhu et al. (2024); Nika et al. (2024) we again show this for non-adversarial MDPs, as well as less simplified/constraint settings

Advancements in Addressing Distribution Shifts Several approaches have been proposed to address the issue of out-of-distribution robustness in reward learning, such as ensembles of conservative reward models (Coste et al., 2023), averaging weights of multiple reward models (Ramé et al., 2024), iteratively updating training labels (Zhu et al., 2024), on-policy reward learning (Lang et al., 2024a), and distributionally robust planning (Zhan et al., 2023).

Our work further emphasizes the usefulness of exploring additional assumptions or methods to mitigate the perils of distribution shift, as we show that without any additional assumptions, there are next to no guarantees. We therefore hope that our work can serve as a theoretical baseline, that people can use to express and analyze their new assumptions or methods.

In classical machine learning, research in out-of-distribution generalization has a long history, and a rich literature of methods exists (Li et al., 2022; Zhou et al., 2022; Wang et al., 2022; Liu et al., 2021; Li et al., 2023; Yoon et al., 2023). These methods could potentially be adapted to address distribution shift challenges in reinforcement learning.

Contextual Bandits In Section 6 we work in the contextual bandit setting and derive variants of our results for RLHF. Several theoretical results have been developed that investigate the challenge of RLHF (Xiong et al., 2024; Zhu et al., 2023; Ji et al., 2023; Mehta et al., 2023) and reward learning in general. (Agarwal et al., 2012; Foster et al., 2020) in the contextual bandit setting. Compared to this prior work, we focus on the offline setting where the data distribution D has been pre-generated by a reference policy.

B INTRODUCTION

B.1 PRELIMINARIES

A *Markov Decision Process* (MDP) is a tuple $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ where \mathcal{S} is a set of *states*, \mathcal{A} is a set of *actions*, $\tau : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{A})$ is a *transition function*, $\mu_0 \in \Delta(\mathcal{S})$ is an *initial state distribution*, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a *reward function*, and $\gamma \in (0, 1)$ is a *discount rate*. A *policy* is a function $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$. A *trajectory* $\xi = \langle s_0, a_0, s_1, a_1, \dots \rangle$ is a possible path in an MDP. The *return function* G gives the cumulative discounted reward of a trajectory, $G(\xi) = \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1})$, and the *evaluation function* J gives the expected trajectory return given a policy, $J(\pi) = \mathbb{E}_{\xi \sim \pi} [G(\xi)]$. A policy maximizing J is an *optimal policy*. The *state-action occupancy measure* is a function $\eta : \Pi \rightarrow \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$ which assigns each policy $\pi \in \Pi$ a vector of occupancy measure describing the discounted frequency that a policy takes each action in each state. Formally, $\eta(\pi)(s, a) = \eta^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t \cdot P(s_t = s, a_t = a \mid \xi \sim \pi)$. Note that by writing the reward function R as a vector $\vec{R} \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$, we can split J into a linear function of π : $J(\pi) = \eta^\pi \cdot \vec{R}$. The *value function* V of a policy encodes the expected future discounted reward from each state when following that policy. We use \mathcal{R} to refer to the set of all reward functions. When talking about multiple rewards, we give each reward a subscript R_i , and use J_i , G_i , and V_i^π , to denote R_i 's evaluation function, return function, and π -value function.

B.2 PROBLEM FORMALIZATION

The standard RL process using reward learning works roughly like this:

1. You are given a dataset of transition-reward tuples $\{(s_i, a_i, r_i)\}_{i=0}^n$. Here, each $(s_i, a_i) \in \mathcal{S} \times \mathcal{A}$ is a transition from some (not necessarily known) MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ that has been sampled using some distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$, and $r_i = R(s_i, a_i)$. The goal of the process is to find a policy $\hat{\pi}$ which performs roughly optimally for the unknown true reward function R . More formally: $J_R(\hat{\pi}) \approx \max_{\pi \in \Pi} J_R(\pi)$.
2. Given some error tolerance $\epsilon \in \mathbb{R}$, a reward model $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is learned using the provided dataset. At the end of the learning process \hat{R} satisfies some optimality criterion such as: $\mathbb{E}_{(s,a) \sim D} [|\hat{R}(s, a) - R(s, a)|] < \epsilon$
3. The learned reward model \hat{R} is used to train a policy $\hat{\pi}$ that fulfills the following optimality criterion: $\hat{\pi} = \arg \max_{\pi \in \Pi} J_{\hat{R}}(\pi)$.

The problem is that training $\hat{\pi}$ to optimize \hat{R} effectively leads to a distribution shift, as the transitions are no longer sampled from the original data distribution D but some other distribution \hat{D} (induced by the policy $\hat{\pi}$). Depending on the definition of D , this could mean that there are

no guarantees about how close the expected error of \hat{R} to the true reward function R is (i.e., $\mathbb{E}_{(s,a) \sim \hat{D}} [|\hat{R}(s,a) - R(s,a)|]$ could not be upper-bounded).

This means that we have no guarantee about the performance of $\hat{\pi}$ with respect to the original reward function R , so it might happen that $\hat{\pi}$ performs arbitrarily bad under the true reward R : $J_R(\hat{\pi}) \ll \max_{\pi} J_R(\pi)$.

If for a given data distribution D there exists a reward model \hat{R} such that \hat{R} is close in expectation to the true reward function R but it is possible to learn a policy that performs badly under J_R despite being optimal for \hat{R} , we say that D allows for error-regret mismatch and that \hat{R} has an error-regret mismatch.

C EXISTENCE OF ERROR-REGRET MISMATCH

In this section, we answer the question under which circumstances error-regret mismatch could occur. We consider multiple different settings, starting from very weak statements, and then steadily increasing the strength and generality.

C.1 ASSUMPTIONS

For every MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ that we will define in the following statements, we assume the following properties:

- **Finiteness:** Both the set of states \mathcal{S} and the set of actions \mathcal{A} are finite
- **Reachability:** Every state in the given MDP’s is reachable, i.e., for every state $s \in \mathcal{S}$, there exists a path of transitions from some initial state s_0 (s.t. $\mu_0(s_0) > 0$) to s , such that every transition (s, a, s') in this path has a non-zero probability, i.e., $\tau(s'|s, a) > 0$. Note that this doesn’t exclude the possibility of some transitions having zero probability in general.

C.2 INTUITIVE UNREGULARIZED EXISTENCE STATEMENT

Definition C.1 (Regret). We define the *regret* of a policy π with respect to reward function R as

$$\text{Reg}^R(\pi) := \frac{\max J_R - J_R(\pi)}{\max J_R - \min J_R} \in [0, 1].$$

Here, J is the policy evaluation function corresponding to R .

Definition C.2 (Policy-Induced Distribution). Let π be a policy. Then we define the *policy-induced distribution* D^π by

$$D^\pi := (1 - \gamma) \cdot \eta^\pi.$$

Definition C.3 (Range of Reward Function). Let R be a reward function. Its *range* is defined as

$$\text{range } R := \max R - \min R.$$

Lemma C.4. for any policy π , D^π is a distribution.

Proof. This is clear. □

Proposition C.5. Let $M = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP, $D \in \Delta(\mathcal{S} \times \mathcal{A})$ a data distribution, and $\epsilon > 0$, $L \in [0, 1]$. Assume there exists a policy $\hat{\pi}$ with the property that $\text{Reg}^R(\hat{\pi}) \geq L$ and $D(\text{supp } D^{\hat{\pi}}) < \epsilon$, where $\text{supp } D^{\hat{\pi}}$ is defined as the set of state-action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$ such that $D^{\hat{\pi}}(s, a) > 0$. In other words, there is a “bad” policy for R that is not very supported by D . Then, D allows for error-regret mismatch to occur, i.e., $D \in \text{unsafe}(R, \epsilon, L)$.

Proof. We will show that whenever there exists a policy $\hat{\pi}$ with the following two properties:

- $\text{Reg}^R(\hat{\pi}) \geq L$;
- $D(\text{supp } D^{\hat{\pi}}) < \epsilon$.

Then there exists a reward function \hat{R} for which $\hat{\pi}$ is optimal, and such that

$$\mathbb{E}_{(s,a) \sim D} \left[\frac{|R(s,a) - \hat{R}(s,a)|}{\text{range } R} \right] \leq \epsilon.$$

Define

$$\hat{R}(s,a) := \begin{cases} R(s,a), & (s,a) \notin \text{supp } D^{\hat{\pi}}; \\ \max R, & \text{else.} \end{cases}$$

Then obviously, $\hat{\pi}$ is optimal for \hat{R} . Furthermore, we obtain

$$\begin{aligned} \mathbb{E}_{(s,a) \sim D} \left[\frac{|R(s,a) - \hat{R}(s,a)|}{\text{range } R} \right] &= \sum_{(s,a)} D(s,a) \frac{|R(s,a) - \hat{R}(s,a)|}{\text{range } R} \\ &= \sum_{(s,a) \in \text{supp } D^{\hat{\pi}}} D(s,a) \frac{\max R - R(s,a)}{\text{range } R} \\ &\leq \sum_{(s,a) \in \text{supp } D^{\hat{\pi}}} D(s,a) \\ &= D(\text{supp } D^{\hat{\pi}}) \\ &\leq \epsilon. \end{aligned}$$

That was to show. \square

Corollary C.6. *Let $M = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP, $\epsilon > 0$, and $L \in [0, 1]$. Assume there exists a set of policies Π_L with:*

- $\text{Reg}^R(\pi) \geq L$ for all $\pi \in \Pi_L$;
- $\text{supp } D^\pi \cap \text{supp } D^{\pi'} = \emptyset$ for all $\pi, \pi' \in \Pi_L$; and
- $|\Pi_L| \geq 1/\epsilon$.

Then $\text{unsafe}(R, \epsilon, L) = \Delta(\mathcal{S} \times \mathcal{A})$, i.e.: all distributions are unsafe.

Proof. Let $D \in \Delta(\mathcal{S} \times \mathcal{A})$. Let $\pi \in \arg \min_{\pi' \in \Pi_L} D(\text{supp } D^{\pi'})$. We obtain

$$|\Pi_L| \cdot D(\text{supp } D^\pi) \leq \sum_{\pi' \in \Pi_L} D(\text{supp } D^{\pi'}) = D \left(\bigcup_{\pi' \in \Pi_L} \text{supp } D^{\pi'} \right) \leq 1,$$

and therefore $D(\text{supp } D^\pi) \leq 1/|\Pi_L| < \epsilon$. The result follows from Proposition 3.3. \square

Proposition C.7. *The assumptions on ϵ in Proposition 3.2 and Proposition 3.3 cannot hold simultaneously.*

Proof. If they would hold simultaneously, we would get:

$$\min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \leq D(\text{supp } D^{\hat{\pi}}) < \epsilon < \frac{1-\gamma}{\sqrt{2}} \cdot \frac{\text{range } J_R}{\text{range } R} \cdot \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \cdot L.$$

Here, the first step is clear, the second step is the assumption from Proposition 3.3, and the third step is the assumption from Proposition 3.2. We now show that this leads to a contradiction.

Dividing by the minimum on both sides, we obtain

$$1 < \frac{L}{\sqrt{2}} \cdot \frac{(1-\gamma)\text{range } J_R}{\text{range } R}. \quad (5)$$

Clearly, we have $L/\sqrt{2} < 1$. We also claim that the second fraction is smaller or equal to 1, which then leads to the desired contradiction. Indeed, let π^* and π_* be an optimal and a worst-case policy, respectively. Then we have

$$\begin{aligned}
(1 - \gamma)\text{range}J_R &= (1 - \gamma)(J_R(\pi^*) - J_R(\pi_*)) \\
&= (1 - \gamma)\eta^{\pi^*} \cdot \vec{R} - (1 - \gamma)\eta^{\pi_*} \cdot \vec{R} \\
&= D^{\pi^*} \cdot \vec{R} - D^{\pi_*} \cdot \vec{R} \\
&= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D^{\pi^*}(s,a)R(s,a) - \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D^{\pi_*}(s,a)R(s,a) \\
&\leq \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} R(s,a) - \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} R(s,a) \\
&= \text{range}R.
\end{aligned}$$

Here, we used the formulation of the policy evaluation function in terms of the occupancy measure η , and then that $1 - \gamma$ is a normalizing factor that transforms the occupancy measure into a distribution. Overall, this means that $(1 - \gamma)\text{range}J_R/\text{range}R \leq 1$, contradicting (5). Consequently, the assumptions of Proposition 3.2 and Proposition 3.3 cannot hold simultaneously. \square

C.3 GENERAL EXISTENCE STATEMENTS

We start by giving some definitions:

Definition C.8 (Minkowski addition). Let A, B be sets of vectors, then the Minkowski addition of A, B is defined as:

$$A + B := \{a + b \mid a \in A, b \in B\}.$$

(Karwowski et al., 2023) showed in their proposition 1, that for every MDP, the corresponding occupancy measure space Ω forms a convex polytope. Furthermore, for each occupancy measure $\eta \in \Omega$ there exists at least one policy π^η such that $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$, $\eta^{\pi^\eta}(s, a) = \eta(s, a)$ (see Theorem 6.9.1, Corollary 6.9.2, and Proposition 6.9.3 of (Puterman, 1994)). In the following proofs, we will refer multiple times to vertices of the occupancy measure space Ω whose corresponding policies have high regret. We formalize this in the following definition:

Definition C.9 (High regret vertices). Given a lower regret bound $L \in [0, 1]$, an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and a corresponding occupancy measure Ω , we define the set of high-regret vertices of Ω , denoted by V_R^L , to be the set of vertices v of Ω for which $\text{Reg}^R(\pi^v) \geq L$

Definition C.10 (Active inequalities). Let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP with corresponding occupancy measure space Ω . For every $\eta \in \Omega$, we define the set of transitions (s, a) for which $\eta(s, a) = 0$ by $\text{zeros}(\eta)$.

Definition C.11 (Normal cone). The normal cone of a convex set $C \subset \mathbb{R}^n$ at point $x \in C$ is defined as:

$$N_C(x) := \{n \in \mathbb{R}^n \mid n^T \cdot (x' - x) \leq 0 \text{ for all } x' \in C\} \quad (6)$$

We first state a theorem from prior work that we will use to prove some lemmas in this section:

Theorem C.12 ((Schlaginhaufen & Kamgarpour, 2023)). Let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, \gamma \rangle$ be an MDP without reward function and denote with Ω its corresponding occupancy measure space. Then, for every reward function R and occupancy measure $\eta \in \Omega$, it holds that:

$$\eta \text{ is optimal for } R \iff R \in N_\Omega(\eta), \quad (7)$$

where the normal cone is equal to:

$$N_\Omega(\eta) = \Phi + \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(\eta)}) \quad (8)$$

where Φ is the linear subspace of potential functions used for reward-shaping, and the addition is defined as the Minkowski addition.

Proof. This is a special case of theorem 4.5 of Schlaginhaufen & Kamgarpour (2023), where we consider the unconstrained- and unregularized RL problem. \square

From the previous lemma, we can derive the following corollary which uses the fact that Ω is a closed, and bounded convex polytope (see Proposition 1 of Karwowski et al. (2023)).

Corollary C.13. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and a corresponding occupancy measure space Ω , then for every reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and lower regret bound $L \in [0, 1]$, the following two statements are equivalent:*

- a) *There exists an optimal policy $\hat{\pi}$ for \hat{R} such that $\hat{\pi}$ has regret at least L w.r.t. the original reward function, i.e., $\text{Reg}^R(\hat{\pi}) \geq L$.*
- b) *$\hat{R} \in \Phi + \bigcup_{v \in V_R^L} \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})$, where Φ is the linear subspace of potential functions used for reward-shaping, the addition is defined as the Minkowski addition.*

Proof. Let \hat{R} be chosen arbitrarily. Statement a) can be formally expressed as:

$$\exists \hat{\pi} \in \Pi, \text{Reg}^{\hat{R}}(\hat{\pi}) = 0 \wedge \text{Reg}^R(\hat{\pi}) \geq L.$$

Using Theorem C.12, it follows that:

$$\begin{aligned} \exists \hat{\pi} \in \Pi, \text{Reg}^{\hat{R}}(\hat{\pi}) = 0 \wedge \text{Reg}^R(\hat{\pi}) \geq L \\ \iff \exists \hat{\pi} \in \Pi, \hat{R} \in N_{\Omega}(\eta^{\hat{\pi}}) \wedge \text{Reg}^R(\hat{\pi}) \geq L \\ \iff \hat{R} \in \bigcup_{\eta: \text{Reg}^R(\pi^{\eta}) \geq L} N_{\Omega}(\eta). \end{aligned}$$

It remains to be shown that the union in the previous derivation is equivalent to a union over just all V_R^L . First, note that by definition of the set of high-regret vertices V_R^L (see Definition C.9), it trivially holds that:

$$\bigcup_{v \in V_R^L} N_{\Omega}(v) \subseteq \bigcup_{\eta: \text{Reg}^R(\pi^{\eta}) \geq L} N_{\Omega}(\eta), \quad (9)$$

Next, because Ω is a convex polytope, it can be defined as the intersection of a set of defining half-spaces which are defined by linear inequalities:

$$\Omega = \{\eta \mid a_i^T \cdot \eta \leq b_i, \text{ for } i = 1, \dots, m\}.$$

By defining the active index set of a point $\eta \in \Omega$ as $I_{\Omega}(\eta) = \{a_i \mid a_i^T \cdot \eta = b_i\}$, Rockafellar & Wets (2009) then show that:

$$N_{\Omega}(\eta) = \left\{ y_1 \cdot a_1 + \dots + y_m \cdot a_m \mid y_i \geq 0 \text{ for } i \in I_{\Omega}(\eta), y_i = 0 \text{ for } i \notin I_{\Omega}(\eta) \right\}, \quad (10)$$

(see their theorem 6.46). Note that, because Ω lies in an $|\mathcal{S}| \cdot (|\mathcal{A}| - 1)$ dimensional affine subspace (see Proposition 1 of (Karwowski et al., 2023)), a subset of the linear inequalities which define Ω must always hold with equality, namely, the inequalities that correspond to half-spaces which define the affine subspace in which Ω resides. Therefore, the corresponding active index set, let's denote it by $I_{\Omega, \Phi}(\eta)$ because the subspace orthogonal to the affine subspace in which Ω lies corresponds exactly to Φ , is always non-empty and the same for every $\eta \in \Omega$.

Now, from Equation (10), it follows that for every $\eta \in \Omega$, there exists a vertex v of Ω , such that $N_{\Omega}(\eta) \subseteq N_{\Omega}(v)$. We take this one step further and show that for every η with $\text{Reg}^R(\pi^{\eta}) \geq L$, there must exist a vertex v with $\text{Reg}^R(\pi^v) \geq L$ such that $N_{\Omega}(\eta) \subseteq N_{\Omega}(v)$. We prove this via case distinction on η .

- η is in the interior of Ω . In this case, the index set $I_{\Omega}(\eta)$ reduces to $I_{\Omega, \Phi}(\eta)$ and because we have $I_{\Omega, \Phi}(\eta) \subseteq I_{\Omega}(\eta)$ for every $\eta \in \Omega$, the claim is trivially true.
- η itself is already a vertex in which case the claim is trivially true.
- η is on the boundary of Ω . In this case η can be expressed as the convex combination of some vertices V_{η} which lie on the same face of Ω as η . Note that all occupancy measures with regret $\geq L$ must lie on one side of the half-space defined by the equality $R^T \cdot \eta =$

1242 $L \cdot \eta^{\min} + (1 - L) \cdot \eta^{\max}$, where η^{\min} and η^{\max} are worst-case and best-case occupancy
 1243 measures. By our assumption, η also belongs to this side of the half-space. Because η lies in
 1244 the interior of the convex hull of the vertices V_η , at least one $v \in V_\eta$ must therefore also lie
 1245 on this side of the hyperplane and have regret $\geq L$. Because v and η both lie on the same
 1246 face of Ω , we have $I_\Omega(\eta) \subset I_\Omega(v)$ and therefore also $N_\Omega(\eta) \subseteq N_\Omega(v)$.
 1247

1248 Hence, it must also hold that:

$$1249 \bigcup_{\eta: \text{Reg}^R(\pi^\eta) \geq L} N_\Omega(\eta) \subseteq \bigcup_{v \in V_R^L} N_\Omega(v),$$

1250 which, together with Equation (9) proves the claim. \square

1251 The following lemma relates the set of reward functions to the set of probability distributions D

1252 **Lemma C.14.** *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and a second reduced reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow$
 1253 \mathbb{R} , then the following two statements are equivalent:*

- 1254
- 1255 a) *There exists a data distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$ such that $\mathbb{E}_{(s,a) \sim D} [|R(s,a) - \hat{R}(s,a)|] <$
 1256 $\epsilon \cdot \text{range } R$*
- 1257
- 1258 b) *At least one component \hat{R}_i of \hat{R} is "close enough" to R , i.e., it holds that for some transition
 1259 (s, a) : $|R(s, a) - \hat{R}(s, a)| < \epsilon \cdot \text{range } R$.*

1260 *Proof.* We first show the direction $b) \Rightarrow a)$. Assume that $|R(s^*, a^*) - \hat{R}(s^*, a^*)| < \epsilon \cdot \text{range } R$ for
 1261 a given \hat{R} and transition (s^*, a^*) . In that case, we can construct the data distribution D which we
 1262 define as follows:

$$1263 D(s, a) = \begin{cases} p & \text{if } (s, a) \neq (s^*, a^*) \\ 1 - (|\mathcal{S} \times \mathcal{A}| - 1) \cdot p & \text{if } (s, a) = (s^*, a^*) \end{cases}$$

1264 where we choose $p < \min \left(\frac{\epsilon \cdot \text{range } R - |R(s^*, a^*) - \hat{R}(s^*, a^*)|}{\sum_{(s,a) \neq (s^*, a^*)} |R(s,a) - \hat{R}(s,a)|}, \frac{1}{|\mathcal{S} \times \mathcal{A}|} \right)$. From this it can be easily seen
 1265 that:

$$1266 \begin{aligned} 1267 & \mathbb{E}_{(s,a) \sim D} [|R(s, a) - \hat{R}(s, a)|] \\ 1268 &= (1 - (|\mathcal{S} \times \mathcal{A}| - 1) \cdot p) \cdot |R(s^*, a^*) - \hat{R}(s^*, a^*)| \\ 1269 &+ p \cdot \sum_{(s,a) \neq (s^*, a^*)} |R(s, a) - \hat{R}(s, a)| \\ 1270 &< \epsilon \cdot \text{range } R \end{aligned}$$

1271 We now show the direction $a) \Rightarrow b)$ via contrapositive. Whenever it holds that $|R(s, a) - \hat{R}(s, a)| \geq$
 1272 $\epsilon \cdot \text{range } R$ for all transitions $(s, a) \in \mathcal{S} \times \mathcal{A}$, then the expected difference under an arbitrary data
 1273 distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$ can be lower bounded as follows:

$$1274 \begin{aligned} 1275 & \mathbb{E}_{(s,a) \sim D} [|R(s, a) - \hat{R}(s, a)|] \\ 1276 &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a) \cdot |R(s, a) - \hat{R}(s, a)| \\ 1277 &\geq \epsilon \cdot \text{range } R \cdot \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a) \\ 1278 &= \epsilon \cdot \text{range } R \end{aligned}$$

1279 Because this holds for all possible data distributions D we have $\neg b) \Rightarrow \neg a)$ which proves the
 1280 result. \square

Corollary C.13 describes the set of reward functions \hat{R} for which there exists an optimal policy $\hat{\pi}$ that achieves worst-case regret under the true reward function R . Lemma C.14 on the other hand, describes the set of reward functions \hat{R} , for which there exists a data distribution D such that \hat{R} is close to the true reward function R under D . We would like to take the intersection of those two sets of reward functions, and then derive the set of data distributions D corresponding to this intersection. Toward this goal we first present the following lemma:

Lemma C.15. *For all $\epsilon > 0$, $L \in [0, 1]$, MDP $M = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and all data distributions $D \in \Delta(\mathcal{S} \times \mathcal{A})$, there exists a system of linear inequalities, such that $D \in \text{unsafe}(R, \epsilon, L)$ if and only if the system of linear inequalities is solvable.*

More precisely, let V_R^L be the set of high-regret vertices defined as in Definition C.9. Then, there exists a matrix C , as well as a matrix $U(v)$ and a vector $b(v)$ for every $v \in V_R^L$ such that the following two statements are equivalent:

1. $D \in \text{unsafe}(R, \epsilon, L)$, i.e., there exists a reward function \hat{R} and a policy $\hat{\pi}$ such that:

- (a) $\mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s,a) - R(s,a)|}{\text{range } R} \right] \leq \epsilon;$
- (b) $\text{Reg}^R(\hat{\pi}) \geq L$
- (c) $\text{Reg}^{\hat{R}}(\hat{\pi}) = 0$

2. There exists a vertex $v \in V_R^L$ such that the linear system

$$\begin{bmatrix} U(v) \\ C \cdot \text{diag}(D) \end{bmatrix} \cdot B \leq \begin{bmatrix} b(v) \\ \epsilon \cdot \text{range } R \cdot \mathbf{1} \end{bmatrix} \quad (11)$$

has a solution B . Here, we use the vector notation of the data distribution D .

Proof. We can express any reward function \hat{R} as $\hat{R} = R + B$, i.e. describing \hat{R} as a deviation $B : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ from the true reward function. Note that in this case, we get $\hat{R} - R = B$. Next, note that the expression:

$$\mathbb{E}_{(s,a) \sim D} [|B(s, a)|] \leq \epsilon \cdot \text{range } R \quad (12)$$

describes a “weighted L^1 ball” around the origin in which B must lie:

$$\mathbb{E}_{(s,a) \sim D} [|B(s, a)|] \leq \epsilon \cdot \text{range } R \quad (13)$$

$$\iff \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a) \cdot |B(s, a)| \leq \epsilon \cdot \text{range } R \quad (14)$$

$$\iff B \in \mathcal{C}(D) := \left\{ x \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|} \mid \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a) \cdot |x_{s,a}| \leq \epsilon \cdot \text{range } R \right\}. \quad (15)$$

This “weighted L^1 ball” is a polyhedral set, which can be described by the following set of inequalities:

$$\begin{aligned} D(s_1, a_1) \cdot B(s_1, a_1) + D(s_1, a_2) \cdot B(s_1, a_2) + \dots &\leq \epsilon \cdot \text{range } R \\ -D(s_1, a_1) \cdot B(s_1, a_1) + D(s_1, a_2) \cdot B(s_1, a_2) + \dots &\leq \epsilon \cdot \text{range } R \\ D(s_1, a_1) \cdot B(s_1, a_1) - D(s_1, a_2) \cdot B(s_1, a_2) + \dots &\leq \epsilon \cdot \text{range } R \\ -D(s_1, a_1) \cdot B(s_1, a_1) - D(s_1, a_2) \cdot B(s_1, a_2) + \dots &\leq \epsilon \cdot \text{range } R \\ &\dots \end{aligned}$$

This can be expressed more compactly in matrix form, as:

$$C \cdot \text{diag}(D) \cdot B \leq \epsilon \cdot \text{range } R \cdot \mathbf{1}, \quad (16)$$

where $C \in \mathbb{R}^{2^{|\mathcal{S} \times \mathcal{A}|} \times |\mathcal{S} \times \mathcal{A}|}$, $\text{diag}(D) \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}| \times |\mathcal{S} \times \mathcal{A}|}$, $B \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$, $\mathbf{1} \in \{1\}^{|\mathcal{S} \times \mathcal{A}|}$ and the individual matrices are defined as follows:

$$C = \begin{bmatrix} 1 & 1 & \dots & 1 \\ -1 & 1 & \dots & 1 \\ 1 & -1 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ -1 & -1 & \dots & -1 \end{bmatrix}, \quad \text{diag}(D) = \begin{bmatrix} D(s_1, a_1) & & 0 \\ & \ddots & \\ 0 & & D(s_n, a_m) \end{bmatrix}. \quad (17)$$

Next, from Corollary C.13 we know that a reward function $\hat{R} = R + B$ has an optimal policy with regret larger or equal to L if and only if:

$$\begin{aligned} R + B &\in \Phi + \bigcup_{v \in V_R^L} \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}) \\ \iff B &\in -R + \Phi + \bigcup_{v \in V_R^L} \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}) \end{aligned} \quad (18)$$

We can rephrase the above statement a bit. Let's focus for a moment on just a single vertex $v \in V_R^L$. First, note that because Φ and $\text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})$, are polyhedral, $\Phi + \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})$ must be polyhedral as well (this follows directly from Corollary 3.53 of (Rockafellar & Wets, 2009)). Therefore, the sum on the right-hand side can be expressed by a set of linear constraints $U(v) \cdot B \leq b(v)$.

Hence, a reward function, $\hat{R} = R + B$ is close in expected L1 distance to the true reward function R , and has an optimal policy that has large regret with respect to R , if and only if there exists at least one vertex $v \in V_R^L$, such that:

$$\begin{bmatrix} U(v) \\ C \cdot \text{diag}(D) \end{bmatrix} \cdot B \leq \begin{bmatrix} b(v) \\ \epsilon \cdot \text{range } R \cdot \mathbf{1} \end{bmatrix} \quad (19)$$

holds. □

In the next few subsections, we provide a more interpretable version of the linear system of inequalities in Equation (11), and the conditions for when it is solvable and when not.

C.3.1 MORE INTERPRETABLE STATEMENT

Ideally, we would like to have a more interpretable statement about which classes of data distributions D fulfill the condition of Equation (11). We now show that for an arbitrary MDP and data distribution D , D is a safe distribution, i.e., error-regret mismatch is not possible, if and only if D fulfills a fixed set of linear constraints (independent of D).

Theorem C.16. *For all MDPs $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ and $L \in [0, 1]$, there exists a matrix M such that for all $\epsilon > 0$ and $D \in \Delta(\mathcal{S} \times \mathcal{A})$ we have:*

$$D \in \text{safe}(R, \epsilon, L) \iff M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}, \quad (20)$$

where we use the vector notation of D , and $\mathbf{1}$ is a vector containing all ones.

Proof. Remember from Lemma C.15, that a data distribution D is safe, i.e., $D \in \text{safe}(R, \epsilon, L)$, if and only if for all unsafe vertices $v \in V_R^L$ the following system of linear inequalities:

$$\begin{bmatrix} U(v) \\ C \cdot \text{diag}(D) \end{bmatrix} \cdot B \leq \begin{bmatrix} b(v) \\ \epsilon \cdot \text{range } R \cdot \mathbf{1} \end{bmatrix} \quad (21)$$

has no solution. Let $v \in V_R^L$ be chosen arbitrarily and define $\mathcal{U}_v := \{B \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|} \mid U(v) \cdot B \leq b(v)\}$, i.e., \mathcal{U}_v is the set of all $B \in \mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$, such that $\hat{R} := R + B$ has an optimal policy with regret at least L . Then, Equation (21) has no solution if and only if:

$$\forall B \in \mathcal{U}_v, \quad C \cdot \text{diag}(D) \cdot B \not\leq \epsilon \cdot \text{range } R \cdot \mathbf{1} \quad (22)$$

$$\iff \forall B \in \mathcal{U}_v, \quad \text{abs}(B)^T \cdot D > \epsilon \cdot \text{range } R, \quad (23)$$

where we used the definition of the matrices C , and $\text{diag}(D)$ (see Equation (16)) and $\text{abs}(\cdot)$ denotes the element-wise absolute value function. Now, we will finish the proof by showing that there exists a *finite* set of vectors $X \subset \mathcal{U}_v$ (which is independent of the choice of D), such that for every $x \in X$, Equation (23) holds if and only if it is true for all B , i.e., more formally:

$$\forall B \in X, \quad \text{abs}(B)^T \cdot D > \epsilon \cdot \text{range } R$$

$$\iff \forall B \in \mathcal{U}_v, \quad \text{abs}(B)^T \cdot D > \epsilon \cdot \text{range } R.$$

And since X is finite, we can then summarize the individual elements of X as rows of a matrix M and get the desired statement by combining the previous few statements, namely:

$$D \in \text{safe}(R, \epsilon, L) \iff M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1} \quad (24)$$

Towards this goal, we start by reformulating Equation (23) as a condition on the optimal value of a convex optimization problem:

$$\begin{aligned} & \forall x \in \mathcal{U}_v, \quad \text{abs}(x)^T \cdot D > \epsilon \cdot \text{range } R \\ \iff & \left(\min_{x \in \mathcal{U}_v} \text{abs}(x)^T \cdot D \right) > \epsilon \cdot \text{range } R \\ \iff & \text{abs}(x^*)^T \cdot D > \epsilon \cdot \text{range } R, \quad \text{where } x^* := \arg \min_{x \in \mathcal{U}_v} \text{abs}(x)^T \cdot D \\ \iff & \text{abs}(x^*)^T \cdot D > \epsilon \cdot \text{range } R, \quad \text{where } x^* := \arg \min_x \text{abs}(x)^T \cdot D, \quad (25) \\ & \text{subject to } U(v) \cdot x \leq b(v) \end{aligned}$$

Note that the optimal value x^* of this convex optimization problem depends on the precise definition of the data distribution D . But importantly, the set over which we optimize (i.e., \mathcal{U}_v defined as the set of all x , such that $U(v) \cdot x \leq b$) does *not* depend on D ! The goal of this part of the proof is to show that for all possible D the optimal value of the optimization problem in Equation (25) is *always* going to be one of the vertices of \mathcal{U}_v . Therefore, we can transform the optimization problem in Equation (25) into a new optimization problem that does not depend on D anymore. It will then be possible to transform this new optimization problem into a simple set of linear inequalities which will form the matrix M in Equation (24).

Towards that goal, we continue by splitting up the convex optimization problem into a set of linear programming problems. For this, we partition $\mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$ into its different orthants O_c for $c \in \{-1, 1\}^{|\mathcal{S} \times \mathcal{A}|}$ (a high-dimensional generalization of the quadrants). More precisely, for every $x \in O_c$, we have $\text{diag}(c) \cdot x = \text{abs}(x)$. Using this definition, we can reformulate the constraint on the convex optimization problem as follows:

$$\min_{\substack{c \in \{-1, 1\}^{|\mathcal{S} \times \mathcal{A}|} \\ x_c \neq \emptyset}} (\text{diag}(c) \cdot x_c)^T \cdot D > \epsilon \cdot \text{range } R, \quad (26)$$

where the individual x_c are defined as the solution of linear programming problems:

$$\begin{aligned} x_c & := \arg \min_x (\text{diag}(c) \cdot x)^T \cdot D \\ & \text{subject to } U(v) \cdot x \leq b(v) \\ & \quad \text{diag}(c) \cdot x \geq 0, \end{aligned} \quad (27)$$

or $x_c := \emptyset$ in case the linear program is infeasible. Finally, by re-parametrizing each linear program using the variable transform $x' = \text{diag}(c) \cdot x$ we can convert these linear programs into standard form:

$$\begin{aligned} x_c & := \text{diag}(c) \cdot \arg \min_{x'} x'^T \cdot D \\ & \text{subject to } U(v) \cdot \text{diag}(c) \cdot x' \leq b(v) \\ & \quad x' \geq 0, \end{aligned} \quad (28)$$

where we used twice the fact that $\text{diag}(c)^{-1} = \text{diag}(c)$, and hence, $x = \text{diag}(c) \cdot x'$. Because it was possible to transform these linear programming problems described in Equation (27) into standard form using a simple variable transform, we can apply standard linear programming theory to draw the following conclusions (see Theorem 3.4 and Section 6 of Chapter 2 of (Vanderbei, 1998) for reference):

1. The set of constraints in Equations (27) and (28) are either infeasible or they form a polyhedral set of feasible solutions.
2. If the set of constraints in Equations (27) and (28) are feasible, then there exists an optimal feasible solution that corresponds to one of the vertices (also called basic feasible solutions) of the polyhedral constraint sets. This follows from the fact that the objective function is bounded from below by zero.

Let's denote the polyhedral set of feasible solutions defined by the constraints in Equation (27) by $\mathcal{F}_c(v)$. Because $\mathcal{F}_c(v)$ does not depend on the specific choice of the data distribution, this must mean that for every possible data distribution D , we have either $x_c = \emptyset$ or x_c is one of the vertices of $\mathcal{F}_c(v)$, denoted by $\text{vertices}(\mathcal{F}_c(v))$! Note that, by definition of x_c , it holds that:

$$\forall x \in \text{vertices}(\mathcal{F}_c(v)), \quad (\text{diag}(c) \cdot x_c)^T \cdot D \leq (\text{diag}(c) \cdot x)^T \cdot D. \quad (29)$$

Therefore, we can define:

$$X(v) := \bigcup_{c \in \{-1, 1\}^{|\mathcal{S} \times \mathcal{A}|}} \text{vertices}(\mathcal{F}_c(v)) = \{x_1, \dots, x_k\}, \quad \text{and} \quad M_{X(v)} := \begin{bmatrix} \text{abs}(x_1)^T \\ \dots \\ \text{abs}(x_k)^T \end{bmatrix}, \quad (30)$$

where $M_{X(v)}$ contains the element-wise absolute value of all vectors of $X(v)$ as row vectors. Let D be an arbitrary data distribution. Then, we've shown the following equivalences:

$$\begin{aligned} & \forall B \in \mathcal{U}_v, \quad \text{abs}(B)^T \cdot D > \epsilon \cdot \text{range } R && \text{(see Equation (23))} \\ \iff & \min_{\substack{c \in \{-1, 1\}^{|\mathcal{S} \times \mathcal{A}|} \\ x_c \neq \emptyset}} (\text{diag}(c) \cdot x_c)^T \cdot D > \epsilon \cdot \text{range } R && \text{(see Equation (26))} \\ \iff & \min_{x \in X(v)} \text{abs}(x)^T \cdot D > \epsilon \cdot \text{range } R && \text{(due to Equation (29))} \\ \iff & M_X(v) \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1} \end{aligned}$$

Now, by combining the individual sets of vertices $X(v)$, as follows:

$$X := \bigcup_{v \in V_R^L} X(v) = \{x_1, \dots, x_l\}, \quad \text{and} \quad M = \begin{bmatrix} \text{abs}(x_1)^T \\ \dots \\ \text{abs}(x_l)^T \end{bmatrix}, \quad (31)$$

we are now ready to finish the proof by combining all previous steps:

$$\begin{aligned} & D \in \text{safe}(R, \epsilon, L) \\ \iff & \forall v \in V_R^L, \forall B \in \mathcal{U}_v, \quad \text{abs}(B)^T \cdot D > \epsilon \cdot \text{range } R \\ \iff & \forall v \in V_R^L, \quad M_X(v) \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1} \\ \iff & M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}. \end{aligned}$$

That was to show. \square

C.3.2 DERIVING THE CONDITIONS ON D

In Theorem C.16 we've shown that there exists a set of linear constraints $M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}$, such that whenever a data distribution D satisfies these constraints, it is safe. In this subsection, we derive closed-form expressions for the individual rows of M to get a general idea about the different factors determining whether an individual data distribution is safe.

In the proof of Theorem C.16, we showed that M has the form:

$$M = \begin{bmatrix} \text{abs}(x_1)^T \\ \vdots \\ \text{abs}(x_l)^T \end{bmatrix},$$

for some set $X = \{x_1, \dots, x_l\}$, where each $x \in X$ belongs to a vertex of the set of linear constraints defined by the following class of system of linear inequalities:

$$\begin{bmatrix} U(v) \\ -\text{diag}(c) \end{bmatrix} \cdot x \leq \begin{bmatrix} b(v) \\ 0 \end{bmatrix} \quad \begin{array}{l} \text{(Corresponds to the set of unsafe reward functions)} \\ \text{(Corresponds to the orthant } O_c) \end{array} \quad (32)$$

for some $v \in V_R^L$ (the set of unsafe vertices of Ω), and some $c \in \{-1, 1\}^{|\mathcal{S} \times \mathcal{A}|}$ (defining the orthant O_c).

To ease the notation in the following paragraphs, we will use the notation \mathcal{U}_v for the polyhedral set of x such that $U(v) \cdot x \leq b(v)$, and $\mathcal{F}_c(v)$ for the set of solutions to the full set of linear inequalities in Equation (32). Furthermore, we will use $n := |\mathcal{S}|$ and $m := |\mathcal{A}|$.

We start by giving a small helper definition.

Definition C.17 (General position, (Stanley, 2024)). Let \mathcal{H} be a set of hyperplanes in \mathbb{R}^n . Then \mathcal{H} is in general position if:

$$\begin{aligned} \{H_1, \dots, H_p\} \subseteq \mathcal{H}, p \leq n &\implies \dim(H_1 \cap \dots \cap H_p) = n - p \\ \{H_1, \dots, H_p\} \subseteq \mathcal{H}, p > n &\implies H_1 \cap \dots \cap H_p = \emptyset \end{aligned}$$

We will use this definition in the next few technical lemmas. First, we claim that each of the vertices of $\mathcal{F}_c(v)$ must lie on the border of the orthant O_c .

Lemma C.18 (Vertices lie on the intersection of the two constraint sets.). *All vertices of the polyhedral set, defined by the system of linear inequalities:*

$$\begin{bmatrix} U(v) \\ -\text{diag}(c) \end{bmatrix} \cdot x \leq \begin{bmatrix} b(v) \\ 0 \end{bmatrix} \quad (33)$$

must satisfy some of the inequalities of $-\text{diag}(c) \cdot x \leq 0$ with equality.

Proof. Let \mathcal{U}_v be the set of solutions of the upper part of the system of linear equations in Equation (33) and O_c be the set of solutions of the lower part of the system of linear equations in Equation (33). The lemma follows from the fact that \mathcal{U}_v can be expressed as follows (see Equation (18) and the subsequent paragraph):

$$\mathcal{U}_v = -R + \Phi + \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}), \quad (34)$$

where Φ is a linear subspace. Hence, for every x that satisfies the constraints $U(v) \cdot x \leq b(v)$, x lies on the interior of the line segment spanned between $x' = x + \phi$, and $x'' = x - \phi$ for some $\phi \in \Phi$, $\phi \neq \mathbf{0}$. Note that every point on this line segment also satisfies the constraints $U(v) \cdot x \leq b(v)$. Therefore, x can only be a vertex if it satisfies some of the additional constraints, provided by the inequalities $-\text{diag}(c) \cdot x \leq 0$, with equality. \square

Consequently, every vertex of $\mathcal{F}_c(v)$ is the intersection of some k -dimensional surface of \mathcal{U}_v and $k > 0$ standard hyperplanes (hyperplanes whose normal vector belongs to the standard basis).

Lemma C.19 (Basis for Φ . (Schlaginhaufen & Kamgarpour, 2023)). *The linear subspace Φ of potential shaping transformations can be defined as:*

$$\Phi = \text{span}(A - \gamma \cdot P),$$

where $A, P \in \mathbb{R}^{(n-m) \times n}$ for $n = |\mathcal{S}|, m = |\mathcal{A}|$ are matrices defined as:

$$A := \begin{bmatrix} \mathbf{1}^m & \mathbf{0}^m & \dots & \mathbf{0}^m \\ \mathbf{0}^m & \mathbf{1}^m & \dots & \mathbf{0}^m \\ \dots & \dots & \ddots & \dots \\ \mathbf{0}^m & \mathbf{0}^m & \dots & \mathbf{1}^m \end{bmatrix}, \quad P := \begin{bmatrix} \text{---} & \tau(\cdot | s_1, a_1) & \text{---} \\ \text{---} & \tau(\cdot | s_1, a_2) & \text{---} \\ \dots & \dots & \dots \\ \text{---} & \tau(\cdot | s_n, a_m) & \text{---} \end{bmatrix},$$

where $\mathbf{0}^m, \mathbf{1}^m$ are column vectors and $\tau(\cdot | s_i, a_j)$ is a row vector of the form $[\tau(s_1 | s_i, a_j), \dots, \tau(s_n | s_i, a_j)]$.

Furthermore, we have $\dim \Phi = n$.

Proof. This has been proven by (Schlaginhaufen & Kamgarpour, 2023) (see their paragraph "Identifiability" of Section 4). The fact that $\dim \Phi = n$ follows from the fact that Φ is the linear space orthogonal to the affine space containing the occupancy measure space Ω , i.e. $\Phi^\perp = L$ where L is the linear subspace parallel to $\text{span}(\Omega)$ (see the paragraph *Convex Reformulation* of Section 3 of (Schlaginhaufen & Kamgarpour, 2023)) and the fact that $\dim \text{span}(\Omega) = n \cdot (m - 1)$ (see Proposition 1 of (Karwowski et al., 2023)). \square

Lemma C.20 (Dimension of \mathcal{U}_v). $\dim \mathcal{U}_v = n \cdot m$.

Proof. Remember that \mathcal{U}_v can be expressed as follows (see Equation (18) and the subsequent paragraph):

$$\mathcal{U}_v = -R + \Phi + \text{cone}(\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}), \quad (35)$$

From Lemma C.19 we know that $\dim \Phi = n$. We will make the argument that:

- 1566 a) $\dim [\text{cone} (\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})] \geq n \cdot (m - 1)$
 1567
 1568 b) There exist exactly $n \cdot (m - 1)$ basis vectors of $\text{cone} (\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})$ such that the
 1569 combined set of these vectors and the basis vectors of Φ is linearly independent.
 1570

1571 From this, it must follow that:

$$1572 \dim [\Phi + \text{cone} (\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})] = \dim [\Phi] + n \cdot (m - 1) = n \cdot m$$

1575 For a), remember that v is a vertex of the occupancy measure space Ω and that each vertex v of Ω
 1576 corresponds to at least one deterministic policy π^v (see Proposition 1 of (Karwowski et al., 2023)).
 1577 And since every deterministic policy is zero for exactly $n \cdot (m - 1)$ transitions, it must follow that v
 1578 is also zero in *at least* $n \cdot (m - 1)$ transitions, since whenever $\pi^v(a|s) = 0$ for some $(s, a) \in \mathcal{S} \times \mathcal{A}$,
 1579 we have:

$$1580 v(s, a) = \sum_{t=0}^{\infty} \gamma^t \cdot P(s_t = s, a_t = a | \pi^v, \tau) = \pi^v(a|s) \cdot \sum_{t=0}^{\infty} \gamma^t \cdot P(s_t = s | \pi^v, \tau) = 0.$$

1583 Therefore, it follows that $\dim [\text{cone} (\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)})] \geq n \cdot (m - 1)$.
 1584

1585 For b), (Puterman, 1994) give necessary and sufficient conditions for a point $x \in \mathbb{R}^{n \cdot m}$ to be part of
 1586 Ω (see the dual linear program in section 6.9.1 and the accompanying explanation), namely:

$$1587 x \in \Omega \iff [(A - \gamma \cdot P)^T \cdot x = \mu_0 \quad \text{and} \quad I \cdot x \geq 0],$$

1588 where I is the identity matrix and we use the vector notation of the initial state distribution μ_0 .
 1589 Because v is a vertex of Ω , it can be described as the intersection of $n \cdot m$ supporting hyperplanes of
 1590 Ω that are in general position. Because $(A - \gamma \cdot P)$ has rank n (see Lemma C.19), this must mean
 1591 that for v at least $n \cdot (m - 1)$ inequalities of the system $I \cdot v \geq 0$ hold with equality and the combined
 1592 set of the corresponding row vectors and the row vectors of $(A - \gamma \cdot P)^T$ is linearly independent (as
 1593 the vectors correspond to the normal vectors of the set of $n \cdot m$ hyperplanes in general position).
 1594

1595 Note that the set of unit vectors that are orthogonal to v is precisely defined by $\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}$,
 1596 since, by definition of $\text{zeros}(v)$ (see Definition C.10), we have

$$1597 \forall x \in \{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}, \quad x^T \cdot v = 0.$$

1598 From this, it must follow that the polyhedral set \mathcal{U}_v , has dimension $n \cdot m$. □
 1600

1601 **Lemma C.21** (Defining the faces of \mathcal{U}_v). *Each k -dimensional face F of \mathcal{U}_v (with $k \geq n$) can be*
 1602 *expressed as:*

$$1603 -R + \Phi + \text{cone} (E_F), \quad \text{where } E_F \subset \{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}, \quad (36)$$

1604 *such that $|E_F| = k - n$ and the combined set of vectors of E_F and the columns of $A - \gamma \cdot P$ is*
 1605 *linearly independent.*

1606 *Proof.* Remember that \mathcal{U}_v can be expressed as follows (see Equation (18) and the subsequent
 1607 paragraph):

$$1608 \mathcal{U}_v = -R + \Phi + \text{cone} (\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}), \quad (37)$$

1609 This means that we can express \mathcal{U}_v as a polyhedral cone, spanned by non-negative combinations of:

- 1610 • The column vectors of the matrix $A - \gamma \cdot P$.
- 1611 • The column vectors of the matrix $-(A - \gamma \cdot P)$. Since Φ is a linear subspace and a cone
 1612 is spanned by only the positive combinations of its set of defining vectors we also have to
 1613 include the negative of this matrix to allow arbitrary linear combinations.
- 1614 • The set of vectors $\{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}$.

Consequently, each face of \mathcal{U}_v of dimension k is spanned by a subset of the vectors that span \mathcal{U}_v and is therefore also a cone of these vectors. Because the face has dimension k , we require exactly k linearly independent vectors, as it's not possible to span a face of dimension k with less than k linearly independent vectors, and every additional linearly independent vector would increase the dimension of the face. Furthermore, since Φ is a linear subspace that is unbounded by definition, it must be part of every face. Therefore, every face of \mathcal{U}_v has a dimension of at least n (the dimension of Φ). \square

Note that the converse of Lemma C.21 doesn't necessarily hold, i.e., not all sets of the form described in Equation (36) are necessarily surfaces of the polyhedral set $U(v) \cdot x \leq b(v)$.

We are now ready to develop closed-form expressions for the vertices of $\mathcal{F}_c(v)$. Note that it is possible for $\mathbf{0} \in \mathbb{R}^{n \cdot m}$ to be a vertex of $\mathcal{F}_c(v)$. But in this case, according to Theorem C.16, this must mean that the linear system of inequalities $M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}$ is infeasible (since M would contain a zero row and all elements on the right-hand side are non-negative), which means that in this case $\text{safe}(R, \epsilon, L) = \emptyset$. We will therefore restrict our analysis to all non-zero vertices of $\mathcal{F}_c(v)$.

Proposition C.22 (Vertices of $\mathcal{F}_c(v)$). *Every vertex v_{FG} of $\mathcal{F}_c(v)$, with $v_{FG} \neq \mathbf{0}$, lies on the intersection of some face F of the polyhedral set \mathcal{U}_v and some face G of the orthant O_c and is defined as follows:*

$$v_{FG} = -R + [A - \gamma \cdot P, E_F] \cdot \left(E_G \cdot [A - \gamma \cdot P, E_F] \right)^{-1} \cdot E_G \cdot R,$$

where E_F, E_G are matrices whose columns contain standard unit vectors, such that:

$$\begin{aligned} F &= -R + \Phi + \text{cone}(E_F), & \text{for } E_F \subset \{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)} \\ G &= \{x \in \mathbb{R}^{n \cdot m} \mid E_G \cdot x = \mathbf{0}\}. \end{aligned}$$

Proof. We start by defining the faces of the orthant O_c . Remember that O_c is the solution set to the system of inequalities $\text{diag}(c) \cdot x \geq 0$. Therefore, each defining hyperplane of O_c is defined by one row i of $\text{diag}(c)$, i.e. $\text{diag}(c)_i \cdot x = 0$. Note that since $c \in \{-1, 1\}^{n \cdot m}$, this is equivalent to the equation $e_i^T \cdot x = 0$ where e_i is either the i 'th standard unit vector or its negative. And because every 1-dimensional face G of O_c is the intersection of l standard hyperplanes $\{e_{i_1}, \dots, e_{i_l}\}$, this must mean that G is defined as the set of solutions to the system of equations $E_G \cdot x = 0$ where E_G is the matrix whose row vectors are the vectors $\{e_{i_1}, \dots, e_{i_l}\}$.

Next, let v_{FG} be an arbitrary non-zero vertex of $\mathcal{F}_c(v)$. As proven in Lemma C.18, every vertex of $\mathcal{F}_c(v)$ must satisfy some of the inequalities $\text{diag}(c) \cdot x \geq 0$ for $c \in \{-1, 1\}^{n \cdot m}$ with equality. This means that v_{FG} must lie on some face G of the orthant O_c . The non-zero property guarantees that not all inequalities of the system of inequalities $\text{diag}(c) \cdot x \geq 0$ are satisfied with equality, i.e. that G is not a vertex. Assume that $k > 0$ inequalities are *not* satisfied with equality. Therefore, G must have dimension k , and $E_G \in \mathbb{R}^{n \cdot m \times k}$.

Since v_{FG} is a vertex of the intersection of the orthant O_c and the polyhedral set \mathcal{U}_v , and it only lies on a k -dimensional face of O_c , it must also lie on a $n \cdot m - k$ dimensional face F of \mathcal{U}_v such that the combined set of hyperplanes defining F and G is in general position. The condition that the combined set of hyperplanes is in general position is necessary, to guarantee that v_{FG} has dimension 0 and is therefore a proper vertex.

From Lemma C.21 we know that F can be expressed as:

$$-R + \Phi + \text{cone}(E_F), \quad \text{where } E_F \subset \{-e_{s,a}\}_{(s,a) \in \text{zeros}(v)}, \quad (38)$$

such that $|E_F| = n \cdot (m - 1) - k$ and the combined set of vectors of E_F and the columns of $A - \gamma \cdot P$ are linearly independent.

Because v_{FG} is part of both, F and G , we can combine all information that we gathered about F and G and deduce that it must hold that:

$$\underbrace{E_G \cdot v_{FG} = 0}_{\text{equivalent to } v_{FG} \in G}, \quad \text{and} \quad \underbrace{\exists x \in \mathbb{R}^{n \cdot m - k}, \quad v_{FG} = -R + [A - \gamma \cdot P, E_F] \cdot x}_{\text{equivalent to } v_{FG} \in F}, \quad (39)$$

where for x in Equation (39) it additionally must hold that $\forall i \in \{n + 1, \dots, n \cdot m - k\}, x_i \geq 0$. This must hold because these last entries of x should form a convex combination of the vectors in E_F (as

1674 F is defined to lie in the cone of E_F , see Equation (38)). We briefly state the following two facts that
 1675 will be used later in the proof:
 1676

- 1677 a) v_{FG} is the only vector in $\mathbb{R}^{n \cdot m}$ that fulfills both conditions in Equation (39). This is because
 1678 we defined F in such a way that the intersection of F and G is a single point. And only
 1679 points in this intersection fulfill both conditions in Equation (39).
 1680
 1681 b) For every non-zero vertex v_{FG} , there can only exist a single x that satisfies the two conditions
 1682 in Equation (39). This follows directly from the assumption that the combined set of vectors
 1683 of E_F and the columns of $A - \gamma \cdot P$ are linearly independent (see Equation (38) and the
 1684 paragraph below).

1685 We can combine the two conditions in Equation (39) to get the following, unified condition that is
 1686 satisfied for every non-zero vertex v_{FG} :

$$1687 \exists x \in \mathbb{R}^{n \cdot m - k}, \quad E_G \cdot \left(-R + [A - \gamma \cdot P, E_F] \cdot x \right) = \mathbf{0}^{n \cdot m - k}, \quad (40)$$

1688 From this, it is easy to compute the precise coordinates of v_{FG} :

$$1689 x = \left(E_G \cdot [A - \gamma \cdot P, E_F] \right)^{-1} \cdot E_G \cdot R \quad (41)$$

$$1690 \implies v_{FG} = -R + [A - \gamma \cdot P, E_F] \cdot \left(E_G \cdot [A - \gamma \cdot P, E_F] \right)^{-1} \cdot E_G \cdot R. \quad (42)$$

1691 We finish the proof by showing that the matrix inverse in Equation (41) always exists for every
 1692 non-zero vertex v_{FG} . Assume, for the sake of contradiction, that the matrix $E_G \cdot [A - \gamma \cdot P, E_F]$
 1693 is not invertible. We will show that in this case, there exists a $z \in \mathbb{R}^{n \cdot m}$ with $z \neq v_{FG}$ such that z
 1694 fulfills both conditions in Equation (39). As we've shown above in fact a) this is not possible, hence
 1695 this is a contradiction.

1700 Assuming that $E_G \cdot [A - \gamma \cdot P, E_F]$ is not invertible, we know from standard linear algebra that in
 1701 that case the kernel of this matrix has a dimension larger than zero. Let y_1, y_2 , be two elements of
 1702 this kernel with $y_1 \neq y_2$.

1703 Earlier in this proof, we showed that for every non-zero vertex v_{FG} , Equation (40) is satisfiable. Let
 1704 x be a solution to Equation (40). From our assumptions, it follows that both $x + y_1$ and $x + y_2$ must
 1705 also be solutions to Equation (40) as:
 1706

$$1707 \forall y \in \{y_1, y_2\}, \quad E_G \cdot \left(-R + [A - \gamma \cdot P, E_F] \cdot (x + y) \right) \\
 1708 = -E_G \cdot R + E_G \cdot [A - \gamma \cdot P, E_F] \cdot (x + y) \\
 1709 = -E_G \cdot R + E_G \cdot [A - \gamma \cdot P, E_F] \cdot x \\
 1710 = E_G \cdot \left(-R + [A - \gamma \cdot P, E_F] \cdot x \right) \\
 1711 = \mathbf{0}^{n \cdot m - k}.$$

1712 And from this, it will follow that both, $x + y_1$ and $x + y_2$ must satisfy both conditions in Equation (39).
 1713 Because $x + y_1 \neq x + y_2$, it must also hold that:
 1714

$$1715 -R + [A - \gamma \cdot P, E_F] \cdot (x + y_1) \neq -R + [A - \gamma \cdot P, E_F] \cdot (x + y_2),$$

1716 see fact b) above for a proof of this. And this would mean that there exists at least one $z \in \mathbb{R}^{n \cdot m}$
 1717 with $z \neq v_{FG}$ such that z fulfills both conditions in Equation (39). But as we have shown in fact a),
 1718 this is not possible. Therefore, the matrix $E_G \cdot [A - \gamma \cdot P, E_F]$ must be invertible for every non-zero
 1719 vertex v_{FG} . \square
 1720
 1721
 1722

1723 We are now ready to provide more specific information about the exact conditions necessary for a
 1724 data distribution D to be safe.

1725 **Corollary C.23** (Vertices of $\mathcal{F}_c(v)$). *For all $\epsilon > 0$, $L \in [0, 1]$ and MDPs $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, there*
 1726 *exists a matrix M such that:*
 1727

$$D \in \text{safe}(R, \epsilon, L) \iff M \cdot D > \epsilon \cdot \text{range } R \cdot \mathbf{1}, \quad (43)$$

Algorithm 1 Computes the set of conditions used to determine the safety of a data distribution.

```

1728 1: function COMPUTEM( $MDP = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle, L \in [0, 1]$ )
1729 2:    $I \leftarrow$  the set of all unit vectors of dimension  $|\mathcal{S} \times \mathcal{A}|$ . Create a fixed ordering of  $\mathcal{S}$  and  $\mathcal{A}$  and
1730   denote each vector of  $I$  by  $e_{(s,a)}$  for a unique tuple  $(s, a) \in \mathcal{S} \times \mathcal{A}$ .
1731
1732 3:   candidates  $\leftarrow []$ 
1733 4:    $\Pi_d \leftarrow$  Set of deterministic policies of  $MDP$ 
1734 5:   for  $\pi \in \{\pi' \in \Pi_d : \text{Reg}^R(\pi') \geq L\}$  do ▷ Create a set of potential row candidates.
1735 6:      $E \leftarrow \{e_{(s,a)} \in I : \pi(a|s) = 0\}$ 
1736 7:     for  $E_F \subset E$  do
1737 8:       for  $subset \subseteq I \setminus E_F, |subset| = |\mathcal{S}|$  do
1738 9:          $E_G \leftarrow E_F \cup subset$ 
1739 10:         $E_F, E_G \leftarrow \text{ColumnMatrix}(E_F), \text{RowMatrix}(E_G)$ 
1740 11:        candidates.append( $(E_F, E_G)$ )
1741
1742 12:   rows  $\leftarrow []$  ▷ Find the valid rows amongst the candidates
1743 13:   for  $(E_F, E_G) \in$  candidates do
1744 14:      $k \leftarrow \text{num\_columns}(E_F)$ 
1745 15:     if  $\text{rank}(E_G \cdot [A - \gamma \cdot P, -E_F]) = n + k$  then
1746 16:        $x \leftarrow (E_G \cdot [A - \gamma \cdot P, -E_F])^{-1} \cdot E_G \cdot R$ 
1747 17:       if  $\forall i \in \{n, n + 1, \dots, n + k\} x_i \geq 0$  then
1748 18:         row  $\leftarrow \text{abs}(-R + [A - \gamma \cdot P, -E_F] \cdot x)^T$ 
1749 19:         rows.append(row)
1750
1751 20:    $M \leftarrow \text{RowMatrix}(\text{rows})$ 
1752 21:   return  $M$ 

```

for all $D \in \Delta(\mathcal{S} \times \mathcal{A})$, where we use the vector notation of D , and $\mathbf{1}$ is a vector containing all ones.

The matrix M is defined as:

$$M = \begin{bmatrix} \text{abs}(x_1)^T \\ \vdots \\ \text{abs}(x_l)^T \end{bmatrix},$$

where an individual row x_i of M can either be all zeros, or

$$x_i = -R + [A - \gamma \cdot P, E_{i1}] \cdot (E_{i2} \cdot [A - \gamma \cdot P, E_{i1}])^{-1} \cdot E_{i2} \cdot R, \quad (44)$$

where E_{i1}, E_{i2} are special matrices whose columns contain standard unit vectors.

Proof. This is a simple combination of Theorem C.16 and Proposition C.22. □

In particular, Equation (44) shows that whether a particular data distribution D is safe or not depends on the true reward function R , as well as the transition distribution τ (encoded by the matrix P).

C.3.3 ALGORITHM TO COMPUTE THE CONDITIONS ON D

The derivations of Appendix C.3.2 can be used to define a simple algorithm that constructs matrix M . An outline of such an algorithm is presented in Algorithm 1. We use the terms RowMatrix and ColumnMatrix to denote functions that take a set of vectors and arrange them as rows/columns of a matrix.

To give a brief explanation of the algorithm:

- Line 4 follows from the definitions of V_R^L , $X(v)$ and X (see Definition C.9 and eqs. (30) and (31)).

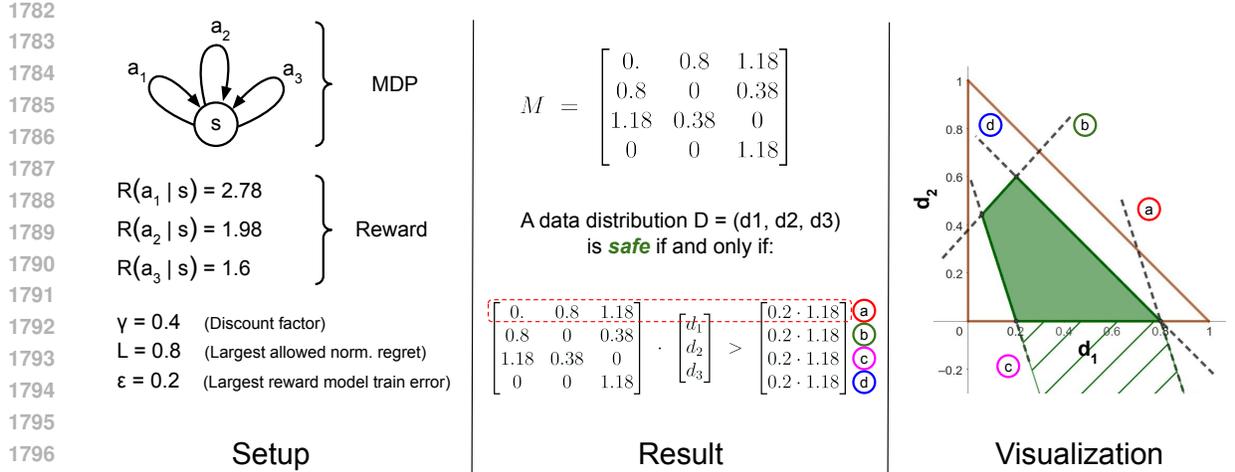


Figure 3: A working example of how to compute the matrix M on a very simple MDP with a single state and three actions. Given the information in the *Setup* column, matrix M can be computed using Algorithm 1. The constructed matrix M contains four linear constraints that a data distribution D has to fulfill in order to be in $\text{safe}(R, \epsilon, L)$. The four constraints are plotted in the right-most column.

- Line 6 are taken from the definition of E_F in Proposition B.20 (except that we don't take the negative of the vectors and instead negate E_F in the final formula).
- Lines 7 and 8 are taken from the definition of E_G (see the first two paragraphs of Proposition C.22). We additionally ensure that E_F is a subset of E_G as otherwise, the matrix $E_G \cdot [A - \gamma \cdot P, -E_F]$ is not invertible (due to the multiplication of $E_G \cdot E_F$) and we know that the matrix must be invertible for every vertex.
- Lines 15 and 17 compute the row of the matrix M . The formulas are a combination of the definition of the sets $X(v)$, X (see Equations (30) and (31)), the matrix M_X (Equation (31)) and Proposition C.22.
- Line 14 checks whether the matrix $E_G \cdot [A - \gamma \cdot P, -E_F]$ is invertible. This is always the case for the rows of M (see the last few paragraphs of the proof of Proposition C.22) but might not be true for other candidates.
- To explain Line 16, remember that every row of the matrix M corresponds to the element-wise absolute value of a vector that lies on the intersection of two polyhedral sets F , and G (see Proposition C.22). The polyhedral set F is defined via a convex cone. To check that our solution candidate lies in this convex cone, we have to check whether the last entries of $x = (E_G \cdot [A - \gamma \cdot P, -E_F])^{-1} \cdot E_G \cdot R$, the entries belonging to the vectors in E_F , are non-negative.

The asymptotic runtime of this naive algorithm is exponential in $|\mathcal{S} \times \mathcal{A}|$ due to the iterations over all subsets in Lines 6 and 7. However, better algorithms might exist and we consider this an interesting direction for future work.

C.3.4 WORKING EXAMPLE OF COMPUTING MATRIX M

Figure 3 shows a simple toy-MDP with a single state and three actions, for which we then compute matrix M using Algorithm 1. Due to the simple structure of the MDP, the auxiliary matrix A and the state-transition matrix P (both used in Algorithm 1) become trivial:

$$A = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \text{ and } P = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

The resulting four constraints that a given data distribution over the state-action space of this MDP has to fulfill to be in $\text{safe}(R, \epsilon, L)$ are then visualized in the right-most column of Figure 3. Note that

the constraints are over three-dimensional vectors. However, because D is a probability distribution, it must live in a two-dimensional subspace of this three-dimensional space, and using the identity $d_3 = 1 - d_1 - d_2$ we can transform the constraints as follows:

$$\begin{bmatrix} | & | & | \\ m_1 & m_2 & m_3 \\ | & | & | \end{bmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix} > \begin{bmatrix} | \\ b \\ | \end{bmatrix} \iff \begin{bmatrix} | & | & | \\ m_1 - m_3 & m_2 - m_3 & | \\ | & | & | \end{bmatrix} \cdot \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} > \begin{bmatrix} | \\ b - m_3 \\ | \end{bmatrix}$$

The brown triangle in Figure 3 depicts the 2d-probability simplex of all distributions over the three actions of the MDP.

Note that constraint ① is a redundant constraint that is already covered by the constraint ④ and the border of the simplex. It would therefore be possible to disregard the computation of such constraints entirely, which could speed up the execution of Algorithm 1. In the next section, we discuss this possibility, as well as other potential directions in which we can extend Theorem 3.5.

C.3.5 BUILDING UP ON THEOREM 3.5

There are multiple ways how future work can build up on the results of Theorem 3.5:

Finding sufficient conditions for safety that require less information about the true reward function: It would be very interesting to investigate whether there exists some subset of the set of safe data distributions for which it is possible to more easily determine membership. This could be helpful in practice, as knowing that a provided data distribution is safe directly yields safety guarantees for the resulting optimal policy.

Developing faster methods to construct M : While the algorithm we provide above runs in exponential time it is unclear whether this has to be the case. The set of vectors that are computed by our algorithm is redundant in the sense that some elements can be dropped as the conditions they encode are already covered by other rows of M . Depending on what fraction of computed elements are redundant it might be possible to develop an algorithm that prevents the computation of redundant rows and can therefore drastically reduce computation time. Alternatively, it would be interesting to develop fast algorithms to compute only parts of M . This could be especially interesting to quickly prove the unsafety of a data distribution, which only requires that a single constraint is violated.

Extending Theorem 3.5 to the regularized policy optimization case: This would allow one to extend the use case we described above to an even wider variety of reward learning algorithms, such as RLHF.

A theoretical baseline (a broader view on the previous point): Most of the options above reveal the properties of the “baseline algorithm” of reinforcement learning under unknown rewards: First, a reward model is trained, and second, a policy is optimized against the trained reward model. The matrix M is valid for the simplest such baseline algorithms without any regularization in either the reward model or the policy. As we mentioned in comments to other reviewers, it would be valuable to study other training schemes (e.g., regularized reward modeling, or switching back and forth between policy optimization and reward modeling on an updated data distribution), for which the set of safe data distributions (or “safe starting conditions”) is likely more favorable than for the baseline case. Then, similar to how empirical work compares new algorithms empirically against baseline algorithms, we hope our work can be a basis to theoretically study improved RL algorithms under unknown rewards, e.g. by deriving a more favorable analog of the matrix M and comparing it with our work.

C.4 EXISTENCE OF NEGATIVE RESULTS IN THE RLHF SETTING

C.4.1 GENERALIZATION OF THE ERROR MEASUREMENT

In this subsection we test the extent to which the results of the previous section generalize to different distance definitions. To ensure compatibility with the positive results of Appendix D.3, we consider MDPs with finite time horizon T . In this setting, trajectories are defined as a finite list of states and actions: $\xi = s_0, a_0, s_1, \dots, a_{T-1}$. Let Ξ be the set of all trajectories of length T . As in the previous

sections, $G : \Xi \rightarrow \mathbb{R}$ denotes the trajectory return function, defined as:

$$G(\xi) = \sum_{t=0}^{T-1} \gamma^t \cdot R(s_t, a_t)$$

Proposition C.24. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and a second reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected difference in trajectory evaluation as follows:*

$$\mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|] \leq \frac{1 - \gamma^T}{1 - \gamma} \cdot \mathbb{E}_{(s,a) \sim D^\pi} [|R(s, a) - \hat{R}(s, a)|] \quad (45)$$

where $D^\pi = \frac{1 - \gamma^T}{1 - \gamma} \cdot \eta^\pi$.

Proof. This follows from the subsequent derivation:

$$\begin{aligned} \mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|] &= \sum_{\xi \in \Xi} P(\xi | \pi) \cdot \left| \sum_{t=0}^{T-1} \gamma^t \cdot (R(s_t, a_t) - \hat{R}(s_t, a_t)) \right| \\ &\leq \sum_{\xi \in \Xi} P(\xi | \pi) \cdot \sum_{t=0}^{T-1} \gamma^t \cdot |R(s_t, a_t) - \hat{R}(s_t, a_t)| \\ &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \left(\sum_{t=0}^{T-1} \gamma^t \cdot P(s_t = s, a_t = a | \pi) \right) \cdot |R(s, a) - \hat{R}(s, a)| \\ &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \eta^\pi(s, a) \cdot |R(s, a) - \hat{R}(s, a)| \\ &= \frac{1 - \gamma^T}{1 - \gamma} \cdot \mathbb{E}_{(s,a) \sim D^\pi} [|R(s, a) - \hat{R}(s, a)|] \end{aligned}$$

□

Given some reward function R , define the probability of trajectory ξ_1 being preferred over trajectory ξ_2 to be:

$$p_R(\xi_1 \succ \xi_2) = \sigma(G_R(\xi_1) - G_R(\xi_2)) = \frac{\exp(G_R(\xi_1))}{\exp(G_R(\xi_1)) + \exp(G_R(\xi_2))}.$$

Then, the following statement holds:

Proposition C.25. *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and a second reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected KL divergence over trajectory preference distributions as follows:*

$$\mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [\mathbb{D}_{KL}(p_R(\cdot | \xi_1, \xi_2) || p_{\hat{R}}(\cdot | \xi_1, \xi_2))] \leq 2 \cdot \mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|], \quad (46)$$

Proof. The right-hand-side of Equation (46) can be lower bounded as follows:

$$2 \cdot \mathbb{E}_{\xi \sim \pi} [|G_R(\xi) - G_{\hat{R}}(\xi)|] \quad (47)$$

$$= \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [|G_R(\xi_1) - G_{\hat{R}}(\xi_1)| + |G_R(\xi_2) - G_{\hat{R}}(\xi_2)|] \quad (48)$$

$$\geq \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [|(G_R(\xi_1) - G_R(\xi_2)) - (G_{\hat{R}}(\xi_1) - G_{\hat{R}}(\xi_2))|] \quad (49)$$

$$= \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} [|x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}|], \quad (50)$$

where from Equation (48) to Equation (49) we used the triangle inequality and did some rearranging of the terms, and from Equation (49) to Equation (50) we simplified the notation a bit by defining $x_{\xi_1, \xi_2} := G_R(\xi_1) - G_R(\xi_2)$ and $y_{\xi_1, \xi_2} := G_{\hat{R}}(\xi_1) - G_{\hat{R}}(\xi_2)$.

Similarly, we can reformulate the left-hand-side of Equation (46) as follows:

$$\mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\mathbb{D}_{\text{KL}} \left(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2) \right) \right] \quad (51)$$

$$= \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\sum_{\substack{i, j \in \{1, 2\} \\ i \neq j}} p_R(\xi_i \succ \xi_j | \xi_1, \xi_2) \cdot \log \left(\frac{p_R(\xi_i \succ \xi_j | \xi_1, \xi_2)}{p_{\hat{R}}(\xi_i \succ \xi_j | \xi_1, \xi_2)} \right) \right] \quad (52)$$

$$= \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\sum_{\substack{i, j \in \{1, 2\} \\ i \neq j}} \sigma(G_R(\xi_i) - G_R(\xi_j)) \cdot \log \left(\frac{\sigma(G_R(\xi_i) - G_R(\xi_j))}{\sigma(G_{\hat{R}}(\xi_i) - G_{\hat{R}}(\xi_j))} \right) \right] \quad (53)$$

$$= \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\sum_{\substack{i, j \in \{1, 2\} \\ i \neq j}} \sigma(x_{\xi_i, \xi_j}) \cdot \log \left(\frac{\sigma(x_{\xi_i, \xi_j})}{\sigma(y_{\xi_i, \xi_j})} \right) \right]. \quad (54)$$

We will now prove the lemma by showing that for all $(\xi_1, \xi_2) \in \Xi \times \Xi$ we have:

$$\sum_{\substack{i, j \in \{1, 2\} \\ i \neq j}} \sigma(x_{\xi_i, \xi_j}) \cdot \log \left(\frac{\sigma(x_{\xi_i, \xi_j})}{\sigma(y_{\xi_i, \xi_j})} \right) \leq |x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}|, \quad (55)$$

from which it directly follows that Equation (54) is smaller than Equation (50).

Let $(\xi_1, \xi_2) \in \Xi \times \Xi$ be chosen arbitrarily. We can then upper bound the left-hand side of Equation (55) as follows:

$$\sigma(x_{\xi_1, \xi_2}) \cdot \log \left(\frac{\sigma(x_{\xi_1, \xi_2})}{\sigma(y_{\xi_1, \xi_2})} \right) + \sigma(x_{\xi_2, \xi_1}) \cdot \log \left(\frac{\sigma(x_{\xi_2, \xi_1})}{\sigma(y_{\xi_2, \xi_1})} \right) \quad (56)$$

$$\leq \log \left(\frac{\sigma(x_{\xi_1, \xi_2})}{\sigma(y_{\xi_1, \xi_2})} \right) + \log \left(\frac{\sigma(x_{\xi_2, \xi_1})}{\sigma(y_{\xi_2, \xi_1})} \right) \quad (57)$$

$$= \log \left(\frac{\sigma(x_{\xi_1, \xi_2}) \cdot \sigma(-x_{\xi_1, \xi_2})}{\sigma(y_{\xi_1, \xi_2}) \cdot \sigma(-y_{\xi_1, \xi_2})} \right) \quad (58)$$

$$= \log \left(\frac{\exp(x_{\xi_1, \xi_2}) \cdot (1 + \exp(y_{\xi_1, \xi_2}))^2}{\exp(y_{\xi_1, \xi_2}) \cdot (1 + \exp(x_{\xi_1, \xi_2}))^2} \right) \quad (59)$$

$$= x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2} + 2 \cdot \log \left(\frac{1 + \exp(y_{\xi_1, \xi_2})}{1 + \exp(x_{\xi_1, \xi_2})} \right), \quad (60)$$

where we used the fact that $x_{\xi_1, \xi_2} = G_R(\xi_1) - G_R(\xi_2)$ and therefore, $-x_{\xi_1, \xi_2} = x_{\xi_2, \xi_1}$ (similar for y_{ξ_1, ξ_2}). We now claim that for all $(\xi_1, \xi_2) \in \Xi \times \Xi$ it holds that:

$$x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2} + 2 \cdot \log \left(\frac{1 + \exp(y_{\xi_1, \xi_2})}{1 + \exp(x_{\xi_1, \xi_2})} \right) \leq |x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}| \quad (61)$$

We prove this claim via proof by cases:

$x_{\xi_1, \xi_2} > y_{\xi_1, \xi_2}$: In this case we have $|x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}| = x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}$ and Equation (61) becomes:

$$2 \cdot \log \left(\frac{1 + \exp(y_{\xi_1, \xi_2})}{1 + \exp(x_{\xi_1, \xi_2})} \right) \leq 0.$$

And since $x_{\xi_1, \xi_2} > y_{\xi_1, \xi_2}$ the fraction inside the logarithm is smaller than 1, this equation must hold.

$x_{\xi_1, \xi_2} = y_{\xi_1, \xi_2}$: In this case, Equation (61) reduces to $0 \geq 0$ which is trivially true.

20198 $x_{\xi_1, \xi_2} < y_{\xi_1, \xi_2}$: In this case, we have $|x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2}| = y_{\xi_1, \xi_2} - x_{\xi_1, \xi_2}$ and we can reformulate
 20199 Equation (61) as follows:

$$\begin{aligned}
 20200 & x_{\xi_1, \xi_2} - y_{\xi_1, \xi_2} + 2 \cdot \log \left(\frac{1 + \exp(y_{\xi_1, \xi_2})}{1 + \exp(x_{\xi_1, \xi_2})} \right) \leq y_{\xi_1, \xi_2} - x_{\xi_1, \xi_2} \\
 20201 & \iff \frac{1 + \exp(y_{\xi_1, \xi_2})}{1 + \exp(x_{\xi_1, \xi_2})} \leq \frac{\exp(y_{\xi_1, \xi_2})}{\exp(x_{\xi_1, \xi_2})} \\
 20202 & \iff \exp(x_{\xi_1, \xi_2}) \leq \exp(y_{\xi_1, \xi_2}).
 \end{aligned}$$

20203 Because we assume that $x_{\xi_1, \xi_2} < y_{\xi_1, \xi_2}$, the last equation, and therefore also the first, must be true.

20204 Combining all the previous statements concludes the proof. \square

20205 Finally, in some RLHF scenarios, one prefers to only compare trajectories with a common starting
 20206 state. In the last lemma, we upper-bound the expected error in choice distributions with trajectories
 20207 that share a common starting state by the expected error in choice distributions with arbitrary
 20208 trajectories:

20209 **Proposition C.26.** *Given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, a data sampling policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and
 20210 a second reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we can upper bound the expected KL divergence of
 20211 preference distributions over trajectories with a common starting state as follows:*

$$\begin{aligned}
 20212 & \mathbb{E}_{\substack{s_0 \sim \mu_0, \\ \xi_1, \xi_2 \sim \pi(s_0)}} \left[\mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \right] \leq \frac{1}{\min_{\substack{s' \in \mathcal{S} \\ \mu_0(s') > 0}} \mu_0(s')} \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \right]. \\
 20213 & \hspace{15em} (62)
 \end{aligned}$$

20214 *Proof.* Let $s_0 : \Xi \rightarrow \mathcal{S}$ define the function which outputs the starting state $s \in \mathcal{S}$ of a trajectory
 20215 $\xi \in \Xi$. We can then prove the lemma by directly lower-bounding the right-hand side of Equation (62):

$$\begin{aligned}
 20216 & \mathbb{E}_{\xi_1, \xi_2 \sim \pi \times \pi} \left[\mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \right] \\
 20217 & = \sum_{s_1, s_2 \in \mathcal{S} \times \mathcal{S}} \mu_0(s_1) \cdot \mu_0(s_2) \cdot \sum_{\substack{\xi_1, \xi_2 \in \Xi \times \Xi \\ s_0(\xi_1) = s_1 \\ s_0(\xi_2) = s_2}} p_{\pi, \tau}(\xi_1 | s_1) \cdot p_{\pi, \tau}(\xi_2 | s_2) \cdot \mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \\
 20218 & = \sum_{s_1 = s_2} \mu_0(s_1) \cdot \mu_0(s_2) \cdot \sum_{\substack{\xi_1, \xi_2 \in \Xi \times \Xi \\ s_0(\xi_1) = s_1 \\ s_0(\xi_2) = s_2}} p_{\pi, \tau}(\xi_1 | s_1) \cdot p_{\pi, \tau}(\xi_2 | s_2) \cdot \mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \\
 20219 & + \sum_{s_1 \neq s_2} \mu_0(s_1) \cdot \mu_0(s_2) \cdot \sum_{\substack{\xi_1, \xi_2 \in \Xi \times \Xi \\ s_0(\xi_1) = s_1 \\ s_0(\xi_2) = s_2}} p_{\pi, \tau}(\xi_1 | s_1) \cdot p_{\pi, \tau}(\xi_2 | s_2) \cdot \mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \\
 20220 & \geq \sum_{s_1 = s_2} \mu_0(s_1) \cdot \mu_0(s_2) \cdot \sum_{\substack{\xi_1, \xi_2 \in \Xi \times \Xi \\ s_0(\xi_1) = s_1 \\ s_0(\xi_2) = s_2}} p_{\pi, \tau}(\xi_1 | s_1) \cdot p_{\pi, \tau}(\xi_2 | s_2) \cdot \mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \\
 20221 & \geq \min_{\substack{s' \in \mathcal{S} \\ \mu_0(s') > 0}} \mu_0(s') \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \sum_{\substack{\xi_1, \xi_2 \in \Xi \times \Xi \\ s_0(\xi_1) = s \\ s_0(\xi_2) = s}} p_{\pi, \tau}(\xi_1 | s) \cdot p_{\pi, \tau}(\xi_2 | s) \cdot \mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \\
 20222 & = \min_{\substack{s' \in \mathcal{S} \\ \mu_0(s') > 0}} \mu_0(s') \cdot \mathbb{E}_{\substack{s_0 \sim \mu_0, \\ \xi_1, \xi_2 \sim \pi(s_0)}} \left[\mathbb{D}_{\text{KL}}(p_R(\cdot | \xi_1, \xi_2) \| p_{\hat{R}}(\cdot | \xi_1, \xi_2)) \right],
 \end{aligned}$$

20223 where we used the fact that the KL divergence is always positive. \square

20224 C.4.2 RLHF BANDIT FORMULATION

20225 RLHF, especially in the context of large language models, is usually modeled in a *contextual bandit*
 20226 setting (Ziegler et al., 2019; Stiennon et al., 2020; Bai et al., 2022; Ouyang et al., 2022; Rafailov

et al., 2023)). A *contextual bandit* $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ is defined by a set of states \mathcal{S} , a set of actions \mathcal{A} , a data distribution $\mu_0 \in \Delta(\mathcal{S})$, and a reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The goal is to learn a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ which maximizes the expected return $J(\pi) = \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [R(s, a)]$. In the context of language models, \mathcal{S} is usually called the set of prompts/contexts, and \mathcal{A} the set of responses. We model the human preference distribution over the set of answers \mathcal{A} using the Bradley-Terry model (Bradley & Terry, 1952). Given a prompt $s \in \mathcal{S}$ and two answers $a_1, a_2 \in \mathcal{A}$, then the probability that a human supervisor prefers answer a_1 to answer a_2 is modelled as:

$$p_R(a_1 \succ a_2 | s) = \frac{\exp(R(s, a_1))}{\exp(R(s, a_1)) + \exp(R(s, a_2))}, \quad (63)$$

where $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is assumed to be the true, underlying reward function of the human.

RLHF is usually done with the following steps:

1. **Supervised finetuning:** Train/Fine-tune a language model π_{ref} using supervised training.
2. **Reward learning:** Given a data distribution over prompts $\mu \in \Delta(\mathcal{S})$, use μ and π_{ref} to sample a set of transitions $(s, a_0, a_1) \in \mathcal{S} \times \mathcal{A} \times \mathcal{A}$ where $s \sim \mu$ and $a_0, a_1 \sim \pi_{\text{ref}}(\cdot|s)$. Use this set of transitions to train a reward model \hat{R} which minimizes the following loss:

$$\mathcal{L}_R(\hat{R}) = -\mathbb{E}_{(s, a_0, a_1, c) \sim \mu, \pi_{\text{ref}}, p_R} \left[\log(\sigma(\hat{R}(s, a_c) - \hat{R}(s, a_{1-c}))) \right], \quad (64)$$

where $c \in \{0, 1\}$ and $p(c = 0 | s, a_0, a_1) = p_R(a_0 \succ a_1 | s)$.

3. **RL finetuning:** Use the trained reward model \hat{R} to further finetune the language model π_{ref} using reinforcement learning. Make sure that the new model does not deviate too much from the original model by penalizing the KL divergence between the two models. This can be done by solving the following optimization problem for some $\lambda > 0$:

$$\pi = \arg \max_{\pi} \mathbb{E}_{s \sim \mu, a \sim \pi(\cdot|s)} [\hat{R}(s, a)] - \lambda \cdot \mathbb{D}_{\text{KL}}(\pi(a|s) || \pi_{\text{ref}}(a|s)) \quad (65)$$

C.4.3 SAFE AND UNSAFE DATA DISTRIBUTIONS FOR RLHF

Definition C.27 (Safe- and unsafe data distributions for RLHF). For a given contextual bandit $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$, let $\epsilon > 0$, $L \in [0, 1]$, $\lambda \in [0, \infty)$, and $\pi_{\text{ref}} : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ an arbitrary reference policy. Similarly to Definition 2.1, we define the set of *safe data distributions* $\text{safe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}}))$ for RLHF as all $D \in \Delta(\mathcal{S} \times \mathcal{A})$ such that for all reward functions $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and policies $\hat{\pi} : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that satisfy the following two properties:

1. **Low expected error:** \hat{R} is similar to R in expected choice probabilities under D , i.e.:

$$\mathbb{E}_{(s, a_1, a_2) \sim D} [\mathbb{D}_{\text{KL}}(p_R(\cdot | s, a_1, a_2) || p_{\hat{R}}(\cdot | s, a_2, a_2))] \leq \epsilon \cdot \text{range } R.$$

2. **Optimality:** $\hat{\pi}$ is optimal with respect to \hat{R} , i.e.:

$$\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \cdot \mathbb{D}_{\text{KL}}(\pi(a|s) || \pi_{\text{ref}}(a|s)).$$

we can guarantee that $\hat{\pi}$ has regret smaller than L , i.e.:

3. **Low regret:** $\hat{\pi}$ has a regret smaller than L with respect to R , i.e., $\text{Reg}^R(\hat{\pi}) < L$.

Similarly, we define the set of *unsafe data distributions* to be the complement of $\text{safe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}}))$:

$$\text{unsafe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}})) := \left\{ D \in \Delta(\mathcal{S} \times \mathcal{A}) \mid D \notin \text{safe}^{\text{RLHF}}(R, \epsilon, L, \lambda, \mathbb{D}_{\text{KL}}(\cdot || \pi_{\text{ref}})) \right\}.$$

Note: Property 1 of Definition C.27 is commonly phrased as minimizing (with respect to \hat{R}) the loss $-\mathbb{E}_{(s, a_1, a_2) \sim D, p_R} [\log(\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2)))]$ (which includes p_R , the probability that a_1 is the preferred action over a_2 , in the expectation). Our version of Property 1 is equivalent to this and can be derived from the former by adding the constant (w.r.t. \hat{R}) term $\mathbb{E}_{(s, a_1, a_2) \sim D, p_R} [\log(\sigma(R(s, a_1) - R(s, a_2)))]$.

2106 C.4.4 NEGATIVE RESULTS
2107

2108 A more advanced result can be achieved by restricting the set of possible pre-trained policies π_{ref} . In
2109 the following proofs, we will define $\pi_{R,\lambda}^{\text{rlhf}}$ to be the optimal policy after doing RLHF on π_{ref} with
2110 some reward function R , i.e.,:

2111 **Definition C.28** (RLHF-optimal policy). For any $\lambda \in \mathbb{R}_+$, reward function R and reference policy
2112 π_{ref} , we define the policy maximizing the RLHF objective by:

$$2113 \pi_{R,\lambda}^{\text{rlhf}} = \arg \max_{\pi} \mathbb{E}_{s \sim \mu, a \sim \pi(\cdot|s)} [R(s, a)] - \lambda \cdot \mathbb{D}_{\text{KL}}(\pi(a|s) || \pi_{\text{ref}}(a|s)) \quad (66)$$

2115 $\pi_{R,\lambda}^{\text{rlhf}}$ does have the following analytical definition (see Appendix A.1 of (Rafailov et al., 2023) for a
2116 derivation):

$$2117 \pi_{R,\lambda}^{\text{rlhf}}(a|s) := \frac{\pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{\sum_{a' \in \mathcal{A}} \pi_{\text{ref}}(a'|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a')\right)}. \quad (67)$$

2120 Before stating the next negative result, we prove a small helper lemma which states that doing RLHF
2121 with some reward function R on a policy π_{ref} is guaranteed to improve the policy return concerning
2122 R :

2123 **Lemma C.29.** For any $\lambda \in \mathbb{R}_+$, reward function R and reference policy π_{ref} , it holds that:

$$2124 J_R(\pi_{R,\lambda}^{\text{rlhf}}) \geq J_R(\pi_{\text{ref}}) \quad (68)$$

2127 *Proof.* We have

$$2128 J_R(\pi_{R,\lambda}^{\text{rlhf}}) - \lambda \mathbb{D}_{\text{KL}}(\pi_{R,\lambda}^{\text{rlhf}} || \pi_{\text{ref}}) = J_{\text{KL}}^R(\pi_{R,\lambda}^{\text{rlhf}}, \pi_{\text{ref}})$$

$$2129 \geq J_{\text{KL}}^R(\pi_{\text{ref}}, \pi_{\text{ref}})$$

$$2130 = J_R(\pi_{\text{ref}}).$$

2132 The result follows from the non-negativity of the KL divergence. \square

2134 We begin by proving a helper lemma that we are going to use in subsequent proofs.

2135 **Lemma C.30.** Let $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ be a contextual bandit

2136 Given a lower regret bound $L \in [0, 1)$, we define for every state $s \in \mathcal{S}$ the reward threshold:

$$2137 R_L(s) := (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a),$$

2138 and define $a_s \in \mathcal{A}$ to be an action such that $R(s, a_s) < R_L(s)$.

2142 Let $\pi_{\text{ref}} : \mathcal{S} \rightarrow \mathcal{A}$ be an arbitrary reference policy for which it holds that for every state $s \in \mathcal{S}$ we
2143 have $\pi_{\text{ref}}(a|s) > 0$.

2144 Then, performing KL-regularized policy optimization, starting from $\pi_{\text{ref}} \in \Pi$ and using the reward
2145 function:

$$2146 \hat{R}(s, a) := \begin{cases} R(s, a) & \text{if } a \neq a_s \\ c_s \in \mathbb{R}_+ & \text{if } a = a_s \end{cases}, \quad (69)$$

2149 results in an optimal policy $\hat{\pi}$ such that $\text{Reg}^R(\hat{\pi}) \geq L$, whenever the constants c_s are larger than the
2150 following lower bound:

$$2151 c_s \geq \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right].$$

2155 *Proof.* Denote by $\pi_{\hat{R},\lambda}^{\text{rlhf}}$ the optimal policy for the following KL-regularized optimization problem:

$$2156 \pi_{\hat{R},\lambda}^{\text{rlhf}} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \cdot \mathbb{D}_{\text{KL}}(\pi(a|s) || \pi_{\text{ref}}(a|s)).$$

2158 The closed-form solution for this optimization problem is known (see Definition C.28). Now, we
2159 prove the statement, by assuming the specific definition of \hat{R} (see Equation (69)), as well as that $\pi_{\hat{R},\lambda}^{\text{rlhf}}$

2160 has a regret at least L , and then work backward to derive a necessary lower bound for the individual
 2161 constants c_s .

2162 We start by defining a small helper policy. Let π_{\top} be a deterministic optimal policy for R and π_{\perp} be
 2163 a deterministic worst-case policy for R . We then define $\pi_L(a|s)$ as a convex combination of π_{\top} and
 2164 π_{\perp} :

$$\begin{aligned}
 2167 \pi_L(a|s) &:= (1 - L) \cdot \pi_{\top}(a|s) + L \cdot \pi_{\perp}(a|s) \\
 2168 &= \begin{cases} 1 & \text{if } R(s, a) = \min_{a' \in \mathcal{A}} R(s, a') = \max_{a' \in \mathcal{A}} R(s, a') \\
 2169 & 1 - L & \text{if } R(s, a) = \max_{a' \in \mathcal{A}} R(s, a') \\
 2170 & L & \text{if } R(s, a) = \min_{a' \in \mathcal{A}} R(s, a') \\
 2171 & 0 & \text{Otherwise} \end{cases} \quad (70) \\
 2172 & \\
 2173 &
 \end{aligned}$$

2175 Next, we show that the regret of π_L is L . Let η_{\top} and η_{\perp} be the corresponding occupancy measures
 2176 of π_{\top} and π_{\perp} . Then, we have:

$$\begin{aligned}
 2178 J_R(\pi_L) &= (1 - L) \cdot R^T \cdot \eta_{\top} + L \cdot R^T \cdot \eta_{\perp}, \\
 2179 & \\
 2180 &
 \end{aligned}$$

2182 from which it directly follows that:

$$\begin{aligned}
 2183 \text{Reg}^R(\pi_L) &= \frac{R^T \cdot \eta_{\top} - [(1 - L) \cdot R^T \cdot \eta_{\top} + L \cdot R^T \cdot \eta_{\perp}]}{R^T \cdot \eta_{\top} - R^T \cdot \eta_{\perp}} = L. \\
 2184 & \\
 2185 & \\
 2186 & \\
 2187 & \\
 2188 & \\
 2189 &
 \end{aligned}$$

2190 Now, having defined π_L , we start the main proof. Assume that $\text{Reg}^R(\pi_{\hat{R}, \lambda}^{\text{rlhf}}) \geq L$, which is
 2191 equivalent to $J(\pi_{\hat{R}, \lambda}^{\text{rlhf}}) \leq J(\pi_L)$. By using the definition of the policy evaluation function, we get:

$$\begin{aligned}
 2192 J(\pi_{\hat{R}, \lambda}^{\text{rlhf}}) &\leq J(\pi_L) \\
 2193 &\iff R^T \cdot (\eta^{\pi_{\hat{R}, \lambda}^{\text{rlhf}}} - \eta^{\pi_L}) \leq 0 \\
 2194 &\iff \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} R(s, a) \cdot \mu_0(s) \cdot (\pi_{\hat{R}, \lambda}^{\text{rlhf}}(a|s) - \pi_L(a|s)) \leq 0 \\
 2195 & \\
 2196 & \\
 2197 & \\
 2198 & \\
 2199 & \\
 2200 & \\
 2201 &
 \end{aligned}$$

2202 We will prove the sufficient condition, that for every $s \in \mathcal{S}$, we have:

$$\begin{aligned}
 2203 \sum_{a \in \mathcal{A}} R(s, a) \cdot (\pi_{\hat{R}, \lambda}^{\text{rlhf}}(a|s) - \pi_L(a|s)) &\leq 0 \quad (71) \\
 2204 & \\
 2205 & \\
 2206 & \\
 2207 & \\
 2208 &
 \end{aligned}$$

2209 Before continuing, note that with our definition of π_L (see Equation (70)) we have:

$$\begin{aligned}
 2210 \sum_{a \in \mathcal{A}} R(s, a) \cdot \pi_L(a|s) &= (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a) =: R_L(s). \\
 2211 & \\
 2212 & \\
 2213 &
 \end{aligned}$$

Now, using this fact as well as the definitions of π_L and $\pi_{\hat{R},\lambda}^{\text{rlhf}}$ (see Definition C.28) we prove under which conditions Equation (71) holds:

$$\begin{aligned}
& \sum_{a \in \mathcal{A}} R(s, a) \cdot \left(\pi_{\hat{R},\lambda}^{\text{rlhf}}(a|s) - \pi_L(a|s) \right) \leq 0 \\
\iff & \sum_{a \in \mathcal{A}} R(s, a) \cdot \left[\frac{\pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right)}{\sum_{a' \in \mathcal{A}} \pi_{\text{ref}}(a'|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a')\right)} - \pi_L(a|s) \right] \leq 0 \\
\iff & \sum_{a \in \mathcal{A}} R(s, a) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right) \\
& \leq \left[\sum_{a \in \mathcal{A}} R(s, a) \cdot \pi_L(a|s) \right] \cdot \sum_{a' \in \mathcal{A}} \pi_{\text{ref}}(a'|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a')\right) \\
\iff & \sum_{a \in \mathcal{A}} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right) \leq 0 \\
\iff & \sum_{\substack{a \in \mathcal{A} \\ R(s, a) > R_L(s)}} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right) \\
& \leq \sum_{\substack{a \in \mathcal{A} \\ R(s, a) < R_L(s)}} (R_L(s) - R(s, a)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right)
\end{aligned}$$

Now, according to the assumptions of the lemma, we know that there exists some action a_s for which $R(s, a_s) < R_L(s)$ and $\pi_{\text{ref}}(a_s|s) > 0$. According to our definition of \hat{R} (see Equation (69)), we have $\hat{R}(s, a_s) = c_s$ and $\hat{R}(s, a) = R(s, a)$ for all other actions. We can use this definition to get a lower bound for c_s :

$$\begin{aligned}
& \sum_{\substack{a \in \mathcal{A} \\ R(s, a) > R_L(s)}} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right) \\
& \leq \sum_{\substack{a \in \mathcal{A} \\ R(s, a) < R_L(s)}} (R_L(s) - R(s, a)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a)\right)
\end{aligned} \tag{72}$$

$$\begin{aligned}
\iff & \sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right) \\
& \leq (R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s) \cdot \exp\left(\frac{1}{\lambda} \cdot \hat{R}(s, a_s)\right)
\end{aligned} \tag{73}$$

$$\iff \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right] \leq \hat{R}(s, a_s). \tag{74}$$

□

We can now use this lemma to prove a more general result:

Proposition C.31. *Let $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ be a contextual bandit.*

Given a lower regret bound $L \in [0, 1)$, we define for every state $s \in \mathcal{S}$ the reward threshold:

$$R_L(s) := (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a),$$

Lastly, $\pi_{\text{ref}} : \mathcal{S} \rightarrow \mathcal{A}$ be an arbitrary reference policy for which it holds that for every state $s \in \mathcal{S}$, $\pi_{\text{ref}}(a|s) > 0$ and there exists at least one action $a_s \in \mathcal{A}$ such that:

2268 a) $\pi_{\text{ref}}(a_s|s)$ is small enough, that the following inequality holds:

$$2270 \log \left[\sum_{a \neq a_s} \pi_{\text{ref}}(a|s) \cdot \exp \left(\frac{1}{\lambda} \cdot (R(s, a) - R(s, a_s)) \right) \cdot \frac{R(s, a) - R_L(s)}{R_L(s) - R(s, a_s)} \right] \leq \frac{\epsilon \cdot \text{range } R}{2 \cdot \lambda \cdot \pi_{\text{ref}}(a_s|s)} + \log(\pi_{\text{ref}}(a_s|s))$$

2274 (75)

2275 b) $R(s, a_s) < R_L(s)$

2277 Then, for all $\epsilon > 0$, $\lambda \in [0, \infty)$, data distributions $\mu \in \Delta(S)$, and true reward functions $R : S \times \mathcal{A} \rightarrow \mathbb{R}$, there exists a reward function $\hat{R} : S \times \mathcal{A} \rightarrow \mathbb{R}$, and a policy $\hat{\pi} : S \rightarrow \Delta(\mathcal{A})$ such that:

- 2281 1. $\mathbb{E}_{s, a_1, a_2 \sim \mu, \pi_{\text{ref}}} [\mathbb{D}_{\text{KL}}(p_R(\cdot|s, a_1, a_2) || p_{\hat{R}}(\cdot|s, a_1, a_2))] \leq \epsilon \cdot \text{range } R$
- 2282 2. $\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \cdot \mathbb{D}_{\text{KL}}(\pi(a|s) || \pi_{\text{ref}}(a|s))$
- 2283 3. $\text{Reg}^R(\hat{\pi}) \geq L$,

2287 *Proof.* We will prove the lemma by construction. Namely, we choose:

$$2292 \hat{R}(s, a) := \begin{cases} R(s, a) & \text{if } a \neq a_s \\ c_s \in \mathbb{R}_+ & \text{if } a = a_s \end{cases} \quad (76)$$

2295 where the different c_s are some positive constants defined as follows:

$$2298 \hat{R}(s, a_s) = c_s \geq l_s := \max \left(R(s, a_s), \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp \left(\frac{1}{\lambda} \cdot R(s, a) \right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right] \right).$$

2300 (77)

2301 Furthermore, the closed-form of the optimal policy $\hat{\pi}$ of the KL-regularized optimization problem is known to be $\pi_{\hat{R}, \lambda}^{\text{rlhf}}$ (see Definition C.28). We now claim that this choice of \hat{R} and $\hat{\pi}$ fulfills properties (1) and (3) of the lemma (property (2) is true by assumption).

2304 Property (3) is true because every reference policy π_{ref} and corresponding reward function R that fulfills the conditions of this proposition also fulfills the conditions of Lemma C.30. Hence, we can directly apply Lemma C.30 and get the guarantee that $\text{Reg}^R(\hat{\pi}) \geq L$.

2308 All that remains to be shown, is that condition (1) can be satisfied by using the definition of \hat{R} and the lower bounds in Equation Equation (77). First, note that we can reformulate the expected error definition in condition (1) as follows:

$$2311 \mathbb{E}_{s, a_1, a_2 \sim \mu, \pi_{\text{ref}}} [\mathbb{D}_{\text{KL}}(p_R(\cdot|s, a_1, a_2) || p_{\hat{R}}(\cdot|s, a_1, a_2))] \\ 2312 = \sum_{s \in S} \mu_0(s) \cdot \sum_{a_1, a_2 \in \mathcal{A} \times \mathcal{A}} \pi_{\text{ref}}(a_1|s) \cdot \pi_{\text{ref}}(a_2|s) \cdot \sum_{i, j \in \{1, 2\}} \sigma(R(s, a_i) - R(s, a_j)) \cdot \log \left(\frac{\sigma(R(s, a_i) - R(s, a_j))}{\sigma(\hat{R}(s, a_i) - \hat{R}(s, a_j))} \right) \\ 2314 = 2 \cdot \sum_{s \in S} \mu_0(s) \cdot \sum_{a_1, a_2 \in \mathcal{A} \times \mathcal{A}} \pi_{\text{ref}}(a_1|s) \cdot \pi_{\text{ref}}(a_2|s) \cdot \sigma(R(s, a_1) - R(s, a_2)) \cdot \underbrace{\log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right)}_{=:\mathcal{IS}(a_1, a_2)} \\ 2317 = 2 \cdot \sum_{s \in S} \mu_0(s) \cdot \sum_{a_1, a_2 \in \mathcal{A} \times \mathcal{A}} \pi_{\text{ref}}(a_1|s) \cdot \pi_{\text{ref}}(a_2|s) \cdot \mathcal{IS}(a_1, a_2).$$

Next, note that for every tuple $(a_1, a_2) \in \mathcal{A}$, the sum $\mathcal{IS}(a_1, a_2) + \mathcal{IS}(a_2, a_1)$ can be reformulated as follows:

$$\begin{aligned}
& \mathcal{IS}(a_1, a_2) + \mathcal{IS}(a_2, a_1) \\
&= \sigma(R(s, a_1) - R(s, a_2)) \cdot \log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right) \\
&\quad + \sigma(R(s, a_2) - R(s, a_1)) \cdot \log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \\
&= \sigma(R(s, a_1) - R(s, a_2)) \cdot \log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right) \\
&\quad + \left(1 - \sigma(R(s, a_1) - R(s, a_2)) \right) \cdot \log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \\
&= \sigma(R(s, a_1) - R(s, a_2)) \cdot \underbrace{\left[\log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right) - \log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \right]}_{(A)} \\
&\quad + \underbrace{\log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right)}_{(B)}.
\end{aligned}$$

The term (A) can now be simplified as follows:

$$\begin{aligned}
& \log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right) - \log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \\
&= \log \left(\frac{\sigma(R(s, a_1) - R(s, a_2))}{1 - \sigma(R(s, a_1) - R(s, a_2))} \right) + \log \left(\frac{1 - \sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))}{\sigma(\hat{R}(s, a_1) - \hat{R}(s, a_2))} \right) \\
&= [R(s, a_1) - R(s, a_2)] - [\hat{R}(s, a_1) - \hat{R}(s, a_2)],
\end{aligned}$$

where we used the definition of the inverse of the logistic function. Similarly, the term (B) can be simplified as follows:

$$\begin{aligned}
& \log \left(\frac{\sigma(R(s, a_2) - R(s, a_1))}{\sigma(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \\
&= \log \left(\frac{\exp(R(s, a_2) - R(s, a_1))}{1 + \exp(R(s, a_2) - R(s, a_1))} \cdot \frac{1 + \exp(\hat{R}(s, a_2) - \hat{R}(s, a_1))}{\exp(\hat{R}(s, a_2) - \hat{R}(s, a_1))} \right) \\
&= [R(s, a_2) - R(s, a_1)] - [\hat{R}(s, a_2) - \hat{R}(s, a_1)] + \log \left(\frac{1 + \exp(\hat{R}(s, a_2) - \hat{R}(s, a_1))}{1 + \exp(R(s, a_2) - R(s, a_1))} \right)
\end{aligned}$$

These expressions, together with the fact that $\mathcal{IS}(a, a) = 0$ for all $a \in \mathcal{A}$, allow us to choose an arbitrary ordering \prec on the set of actions \mathcal{A} , and then re-express the sum:

$$\sum_{a_1, a_2 \in \mathcal{A} \times \mathcal{A}} \pi_{\text{ref}}(a_1 | s) \cdot \pi_{\text{ref}}(a_2 | s) \cdot \mathcal{IS}(a_1, a_2) = \sum_{\substack{a_1, a_2 \in \mathcal{A} \times \mathcal{A} \\ a_1 \prec a_2}} \pi_{\text{ref}}(a_1 | s) \cdot \pi_{\text{ref}}(a_2 | s) \cdot (\mathcal{IS}(a_1, a_2) + \mathcal{IS}(a_2, a_1)). \quad (78)$$

2376

Summarizing all the equations above, we get:

2377

2378

2379

2380

2381

2382

2383

2384

2385

2386

2387

2388

2389

2390

2391

Now, by using our particular definition of \hat{R} (see Equation (76)), we notice that whenever both $a_1 \neq a_s$, and $a_2 \neq a_s$, the inner summand of Equation (79) is zero. What remains of Equation (79) can be restated as follows:

2392

2393

2394

2395

2396

2397

2398

2399

2400

2401

2402

2403

2404

2405

2406

2407

2408

2409

2410

2411

2412

2413

2414

2415

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429

To prove property (1), we must show that Equation (80) is smaller or equal to $\epsilon \cdot \text{range } R$. We do this in two steps. First, note that for all states s it holds that $c_s \geq R(s, a_s)$ (this is obvious from the definition of c_s , see Equation (77)). This allows us to simplify Equation (80) by dropping the logarithm term.

$$\begin{aligned}
& \mathbb{E}_{s, a_1, a_2 \sim \mu, \pi_{\text{ref}}} [\mathbb{D}_{\text{KL}}(p_R(\cdot|s, a_1, a_2) || p_{\hat{R}}(\cdot|s, a_1, a_2))] \\
&= 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \left[(R(s, a_s) - c_s) \cdot \left(\sigma(R(s, a_s) - R(s, a)) - 1 \right) \right. \\
&\quad \left. + \log \left(\frac{1 + \exp(R(s, a) - c_s)}{1 + \exp(R(s, a) - R(s, a_s))} \right) \right] \\
&= 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot (c_s - R(s, a_s)) \cdot \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \left(1 - \sigma(R(s, a_s) - R(s, a)) \right) \\
&\quad + 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \log \left(\frac{1 + \exp(R(s, a) - c_s)}{1 + \exp(R(s, a) - R(s, a_s))} \right). \tag{81}
\end{aligned}$$

Now, we choose to define $c_s := l_s + \delta_s$, where l_s is defined in Equation (77) and $\delta_s \geq 0$ such that:

$$\begin{aligned}
& 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot (l_s + \delta_s - R(s, a_s)) \cdot \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \underbrace{\left(1 - \sigma(R(s, a_s) - R(s, a)) \right)}_{< 1} \\
&\quad + 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot \sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \underbrace{\log \left(\frac{1 + \exp(R(s, a) - l_s - \delta_s)}{1 + \exp(R(s, a) - R(s, a_s))} \right)}_{\leq 0 \text{ (because } c_s := l_s + \delta_s \geq R(s, a_s))} \\
&\leq 2 \cdot \sum_{s \in \mathcal{S}} \mu_0(s) \cdot \pi_{\text{ref}}(a_s|s) \cdot (l_s - R(s, a_s)) \stackrel{!}{\leq} \epsilon \cdot \text{range } R. \tag{82}
\end{aligned}$$

Note that the first inequality is always feasible, as we could just choose $\delta_s = 0$ for all $s \in \mathcal{S}$ in which case the inequality must hold due to the last term in the first line being smaller than one and the last

term in the second line being negative. Now, to prove Equation (82), we prove the sufficient condition that for every state $s \in \mathcal{S}$:

$$\pi_{\text{ref}}(a_s|s) \cdot (l_s - R(s, a_s)) \stackrel{!}{\leq} \frac{\epsilon \cdot \text{range } R}{2}. \quad (83)$$

In case that $l_s = R(s, a_s)$, the left-hand side of Equation (83) cancels and the inequality holds trivially. We can therefore focus on the case where $l_s > R(s, a_s)$. In this case, we get:

$$\begin{aligned} & \pi_{\text{ref}}(a_s|s) \cdot \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right] \stackrel{!}{\leq} \frac{\epsilon \cdot \text{range } R}{2} \\ \Leftrightarrow & \log \left[\sum_{a \neq a_s} \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot (R(s, a) - R(s, a_s))\right) \cdot \frac{R(s, a) - R_L(s)}{R_L(s) - R(s, a_s)} \right] \\ & \stackrel{!}{\leq} \frac{\epsilon \cdot \text{range } R}{2 \cdot \lambda \cdot \pi_{\text{ref}}(a_s|s)} + \log(\pi_{\text{ref}}(a_s|s)) \end{aligned}$$

which holds by assumption (a) of the lemma. Therefore, property (1) of the lemma must hold as well which concludes the proof. \square

Proposition C.32. *Let $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ be a contextual bandit.*

Given a lower regret bound $L \in [0, 1)$, we define for every state $s \in \mathcal{S}$ the reward threshold:

$$R_L(s) := (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a),$$

Lastly, let $\pi_{\text{ref}} : \mathcal{S} \rightarrow \mathcal{A}$ be an arbitrary reference policy for which it holds that for every state $s \in \mathcal{S}$, $\pi_{\text{ref}}(a|s) > 0$, and there exists at least one action $a_s \in \mathcal{A}$ such that:

a) $\pi_{\text{ref}}(a_s|s) > 0$, but $\pi_{\text{ref}}(a_s|s)$ is also small enough, that the following inequality holds:

$$\pi_{\text{ref}}(a_s|s) \leq \frac{(R_L(s) - R(s, a_s))}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{4 \cdot \lambda^2} \quad (84)$$

b) $R(s, a_s) < R_L(s)$

Then Π is a subset of the set of policies in Proposition C.31.

Proof. We show this via a direct derivation:

$$\begin{aligned} & \pi_{\text{ref}}(a_s|s) \leq \frac{R_L(s) - R(s, a_s)}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{4 \cdot \lambda^2} \\ \Rightarrow & \frac{1}{\sqrt{\text{range } R}} \cdot \lambda \cdot \sqrt{\frac{\pi_{\text{ref}}(a_s|s) \cdot L \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{R_L(s) - R(s, a_s)}} \leq \frac{\epsilon}{2} \\ \Rightarrow & \pi_{\text{ref}}(a_s|s) \cdot \lambda \cdot \sqrt{\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)}} \leq \frac{\epsilon \cdot \text{range } R}{2} \end{aligned}$$

We continue by lower-bounding the square-root term as follows:

$$\begin{aligned}
& \lambda \cdot \sqrt{\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s)}} \\
& \geq \lambda \cdot \log \left[\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s)} \right] \\
& \geq \lambda \cdot \log \left[\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot [\max_{a \in \mathcal{A}} R(s, a) - R(s, a_s)]\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s)} \right] \\
& \geq \lambda \cdot \log \left[\frac{(\max_{a \in \mathcal{A}} R(s, a) - R_L(s)) \cdot \exp\left(\frac{1}{\lambda} \cdot \max_{a \in \mathcal{A}} R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right] \\
& \geq \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a | s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right]
\end{aligned}$$

By applying this lower bound, we can finish the proof:

$$\begin{aligned}
\pi_{\text{ref}}(a_s | s) & \leq \frac{R_L(s) - R(s, a_s)}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{4 \cdot \lambda^2} \\
\Rightarrow \pi_{\text{ref}}(a_s | s) \cdot \lambda \cdot \sqrt{\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s)}} & \leq \frac{\epsilon \cdot \text{range } R}{2} \\
\Rightarrow \pi_{\text{ref}}(a_s | s) \cdot \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a | s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s | s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right] & \leq \frac{\epsilon \cdot \text{range } R}{2} \\
\Rightarrow \log \left[\sum_{a \neq a_s} \pi_{\text{ref}}(a | s) \cdot \exp\left(\frac{1}{\lambda} \cdot (R(s, a) - R(s, a_s))\right) \cdot \frac{R(s, a) - R_L(s)}{R_L(s) - R(s, a_s)} \right] & \\
& \leq \frac{\epsilon \cdot \text{range } R}{2 \cdot \lambda \cdot \pi_{\text{ref}}(a_s | s)} + \log(\pi_{\text{ref}}(a_s | s))
\end{aligned}$$

That was to show. \square

C.4.5 ANOTHER NEGATIVE RESULT WITH REGULARIZATION

Proposition C.33. *Let $\langle \mathcal{S}, \mathcal{A}, \mu_0, R \rangle$ be a contextual bandit.*

Given a lower regret bound $L \in [0, 1)$, we define for every state $s \in \mathcal{S}$ the reward threshold:

$$R_L(s) := (1 - L) \cdot \max_{a \in \mathcal{A}} R(s, a) + L \cdot \min_{a \in \mathcal{A}} R(s, a),$$

Lastly, let $\pi_{\text{ref}} : \mathcal{S} \rightarrow \mathcal{A}$ be an arbitrary reference policy for which it holds that for every state $s \in \mathcal{S}$, $\pi_{\text{ref}}(a | s) > 0$ and there exists at least one action $a_s \in \mathcal{A}$ such that:

a) $\pi_{\text{ref}}(a_s | s)$ is small enough, that the following inequality holds:

$$\pi_{\text{ref}}(a_s | s) \leq \frac{(R_L(s) - R(s, a_s))}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{\lambda^2} \quad (85)$$

b) $R(s, a_s) < R_L(s)$

Let $D^{\text{ref}}(s, a) := \mu_0(s) \cdot \pi_{\text{ref}}(a | s)$. Then $D^{\text{ref}} \in \text{unsafe}(R, \epsilon, L, \lambda, \omega)$.

Proof. To prove the proposition we show that there exists some reward function \hat{R} , as well as a policy $\hat{\pi}$ such that the following properties hold:

1. $\mathbb{E}_{(s,a) \sim D^{\text{ref}}} \left[\frac{|R(s,a) - \hat{R}(s,a)|}{\text{range } R} \right] \leq \epsilon.$
2. $\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \omega(\pi)$
3. $\text{Reg}^R(\hat{\pi}) \geq L.$

In particular, we choose:

$$\hat{R}(s, a) := \begin{cases} R(s, a) & \text{if } a \neq a_s \\ c_s \in \mathbb{R}_+ & \text{if } a = a_s \end{cases}, \quad (86)$$

where the different c_s are some positive constants defined as follows:

$$\hat{R}(s, a_s) = c_s := \max \left(R(s, a_s), \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right] \right). \quad (87)$$

Furthermore, the closed-form of the optimal policy $\hat{\pi}$ of the KL-regularized optimization problem is known to be $\pi_{\hat{R}, \lambda}^{\text{rlhf}}$ (see Definition C.28). We now claim that this choice of \hat{R} and $\hat{\pi}$ fulfills properties (1) and (3) of the lemma (property (2) is true by assumption).

Property (3) is true because every reference policy π_{ref} and corresponding reward function R that fulfills the conditions of this proposition also fulfills the conditions of Lemma C.30. Hence, we can directly apply Lemma C.30 and get the guarantee that $\text{Reg}^R(\hat{\pi}) \geq L$.

All that remains to be shown, is that condition (1) can be satisfied by using the definition of \hat{R} and in particular, the definition of the individual c_s (see Equation (87)). The expected error expression in condition (1) can be expanded as follows:

$$\mathbb{E}_{(s,a) \sim D^{\text{ref}}} \left[\frac{|R(s, a) - \hat{R}(s, a)|}{\text{range } R} \right] = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mu_0(s) \cdot \pi_{\text{ref}}(a|s) \cdot \frac{|R(s, a) - \hat{R}(s, a)|}{\text{range } R} \stackrel{!}{\leq} \epsilon.$$

We show the sufficient condition that for each state $s \in \mathcal{S}$ it holds:

$$\sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \frac{|R(s, a) - \hat{R}(s, a)|}{\text{range } R} \stackrel{!}{\leq} \epsilon.$$

By using our definition of \hat{R} (see Equation (86)), this further simplifies as follows:

$$\sum_{a \in \mathcal{A}} \pi_{\text{ref}}(a|s) \cdot \frac{|R(s, a) - \hat{R}(s, a)|}{\text{range } R} = \pi_{\text{ref}}(a_s|s) \cdot \frac{\hat{R}(s, a_s) - R(s, a_s)}{\text{range } R} \stackrel{!}{\leq} \epsilon. \quad (88)$$

In the last equation, we were able to drop the absolute value sign because our definition of the constants c_s (see Equation (87)) guarantees that $\hat{R}(s, a_s) \geq R(s, a_s)$.

Next, note that whenever $\hat{R}(s, a_s) = R(s, a_s)$ the left-hand side of Equation (88) cancels out and so the inequality holds trivially. In the following, we will therefore only focus on states where $\hat{R}(s, a_s) > R(s, a_s)$. Note that this allows us to drop the max statement in the definition of the c_s constants (see Equation (87)).

We continue by upper-bounding the difference $\hat{R}(s, a_s) - R(s, a_s)$. By making use of the following identity:

$$R(s, a_s) = \lambda \cdot \log \left[\exp \left(\frac{1}{\lambda} \cdot R(s, a_s) \right) \right],$$

we can move the $R(s, a_s)$ term into the logarithm term of the c_s constants, and thereby upper-bounding the difference $\hat{R}(s, a_s) - R(s, a_s)$ as follows:

$$\begin{aligned}
& \hat{R}(s, a_s) - R(s, a_s) \\
&= \lambda \cdot \log \left[\frac{\sum_{a \neq a_s} (R(s, a) - R_L(s)) \cdot \pi_{\text{ref}}(a|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right] \\
&\leq \lambda \cdot \log \left[\frac{(\max_{a \in \mathcal{A}} R(s, a) - R_L(s)) \cdot \exp\left(\frac{1}{\lambda} \cdot \max_{a \in \mathcal{A}} R(s, a)\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s) \cdot \exp\left(\frac{1}{\lambda} \cdot R(s, a_s)\right)} \right] \\
&\leq \lambda \cdot \log \left[\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot [\max_{a \in \mathcal{A}} R(s, a) - R(s, a_s)]\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right] \\
&\leq \lambda \cdot \log \left[\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)} \right] \\
&\leq \lambda \cdot \sqrt{\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)}}
\end{aligned}$$

We can now put this upper bound back into Equation (88) and convert the inequality into an upper bound for $\pi_{\text{ref}}(a_s|s)$ as follows:

$$\begin{aligned}
& \pi_{\text{ref}}(a_s|s) \cdot \frac{\hat{R}(s, a_s) - R(s, a_s)}{\text{range } R} \\
&\leq \frac{\pi_{\text{ref}}(a_s|s)}{\text{range } R} \cdot \lambda \cdot \sqrt{\frac{L \cdot \text{range } R \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{(R_L(s) - R(s, a_s)) \cdot \pi_{\text{ref}}(a_s|s)}} \\
&= \frac{1}{\sqrt{\text{range } R}} \cdot \lambda \cdot \sqrt{\frac{\pi_{\text{ref}}(a_s|s) \cdot L \cdot \exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)}{R_L(s) - R(s, a_s)}} \stackrel{!}{\leq} \epsilon \\
\implies \pi_{\text{ref}}(a_s|s) &\leq \frac{R_L(s) - R(s, a_s)}{L} \cdot \frac{\text{range } R}{\exp\left(\frac{1}{\lambda} \cdot \text{range } R\right)} \cdot \frac{\epsilon^2}{\lambda^2}.
\end{aligned}$$

The last line in the previous derivation holds by assumption of the proposal. That was to show. \square

C.5 A REGULARIZED NEGATIVE RESULT FOR GENERAL MDPs

Throughout, let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an MDP. Additionally, assume there to be a data distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$ used for learning the reward function. We do a priori *not assume* that D is induced by a reference policy, but we will specialize to that case later on.

We also throughout fix $\epsilon > 0$, $\lambda > 0$, $L \in (0, 1)$, which will represent, respectively, an approximation-error for the reward function, the regularization strength, and a lower regret bound. Furthermore, let $\omega : \Pi \rightarrow \mathbb{R}$ be any continuous regularization function of policies with $\omega(\pi) \geq 0$ for all $\pi \in \Pi$. For example, if there is a nowhere-zero reference policy π_{ref} , then ω could be given by $\omega(\pi) = \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}})$. For any reward function \hat{R} , a policy $\hat{\pi}$ exists that is optimal with respect to regularized maximization of reward:

$$\hat{\pi} \in \arg \max_{\pi} J_{\hat{R}}(\pi) - \lambda \omega(\pi).$$

We will try to answer the following question: Do there exist realistic conditions on ω and D for which there exists \hat{R} together with $\hat{\pi}$ such that the following properties hold?

- $\mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s,a) - R(s,a)|}{\text{range } R} \right] \leq \epsilon.$
- $\text{Reg}^R(\hat{\pi}) \geq L.$

Furthermore, we now fix π_* , a worst-case policy for R , meaning that $\text{Reg}^R(\pi_*) = 1$. We assume π_* to be deterministic.

Lemma C.34. Define $C(L, R) := \frac{(1-L) \cdot \text{range } J_R}{\|R\|}$. Then the following implication holds:

$$\|D^\pi - D^{\pi_*}\| \leq C(L, R) \implies \text{Reg}^R(\pi) \geq L.$$

Proof. Using the Cauchy-Schwarz inequality, the left side of the implication implies:

$$\begin{aligned} J_R(\pi) - \min J_R &= J_R(\pi) - J_R(\pi_*) \\ &= (D^\pi - D^{\pi_*}) \cdot R \\ &\leq \|D^\pi - D^{\pi_*}\| \cdot \|R\| \\ &\leq (1-L) \cdot \text{range } J_R. \end{aligned}$$

By subtracting $\text{range } J_R = \max J_R - \min J_R$ from both sides, then multiplying by -1 , and then dividing by $\text{range } R$, we obtain the result. \square

Lemma C.35. For any (s, a) , we have

$$\frac{D^\pi(s, a)}{1-\gamma} = \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} \tau(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s) \cdot \pi(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s, a),$$

where

$$\tau(s_0, a_0, \dots, s) := \mu_0(s_0) \cdot \left[\prod_{i=1}^{t-1} \tau(s_i | s_{i-1}, a_{i-1}) \right] \cdot \tau(s | s_{t-1}, a_{t-1}),$$

which is the part in the probability of a trajectory that does not depend on the policy, and

$$\pi(s_0, a_0, \dots, s, a) := \pi(a | s) \cdot \prod_{i=0}^{t-1} \pi(a_i | s_i).$$

Proof. We have

$$\begin{aligned} \frac{D^\pi(s, a)}{1-\gamma} &= \sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a | \xi \sim \pi) \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} P(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s, a | \pi) \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} \mu_0(s_0) \pi(a_0 | s_0) \left[\prod_{i=1}^{t-1} \tau(s_i | s_{i-1}, a_{i-1}) \pi(a_i | s_i) \right] \tau(s | s_{t-1}, a_{t-1}) \pi(a | s) \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} \tau(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s) \cdot \pi(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s, a). \end{aligned}$$

\square

Lemma C.36. Let $1 \geq \delta > 0$. Assume that $\pi(a | s) \geq 1 - \delta$ for all $(s, a) \in \text{supp } D^{\pi_*}$ and that π_* is a deterministic policy.² Then for all $(s, a) \in S \times \mathcal{A}$, one has

$$D^{\pi_*}(s, a) - \delta \cdot (1-\gamma) \cdot \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1-\gamma} D^{\pi_*}(s, a) \right) \leq D^\pi(s, a) \leq D^{\pi_*}(s, a) + \frac{\delta}{1-\gamma}. \quad (89)$$

This also results in the following two inequalities:

$$D^\pi(\text{supp } D^{\pi_*}) \geq 1 - \frac{\delta}{1-\gamma}, \quad \|D^\pi - D^{\pi_*}\| \leq \sqrt{|S \times \mathcal{A}|} \cdot \frac{\delta}{1-\gamma}. \quad (90)$$

²In this lemma, one does not need the assumption that π_* is a worst-case policy, but this case will be the only application later on.

2700 *Proof.* Let $(s, a) \in \text{supp } D^{\pi^*}$. We want to apply the summation formula in Lemma C.35, which we
 2701 recommend to recall. For simplicity, in the following we will write s_0, a_0, \dots when we implicitly
 2702 mean trajectories up until s_{t-1}, a_{t-1} . Now, we will write “ π_* -comp” into a sum to indicate that we
 2703 only sum over states and actions that make the whole trajectory-segment *compatible* with policy π_* ,
 2704 meaning all transitions have positive probability and the actions are deterministically selected by π_* .
 2705 Note that if we restrict to such summands, then each consecutive pair $(s_i, a_i) \in \text{supp } D^{\pi^*}$ is in the
 2706 support of D^{π^*} , and thus we can use our assumption $\pi(a_i | s_i) \geq 1 - \delta$ on those. We can use this
 2707 strategy for a lower-bound:

$$\begin{aligned}
 2708 \quad \frac{D^\pi(s, a)}{1 - \gamma} &\geq \sum_{t=0}^{\infty} \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) \cdot \pi(s_0, a_0, \dots, s, a) \\
 2709 &\geq \sum_{t=0}^{\infty} \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) \cdot (1 - \delta)^{t+1} \\
 2710 &\geq \sum_{t=0}^{\infty} \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) \cdot (1 - \delta \cdot (t + 1)).
 \end{aligned} \tag{91}$$

2718 In the last step, we used the classical formula $(1 - \delta)^t \geq 1 - \delta \cdot t$, which can easily be proved by
 2719 induction over t . Now, we split the sum up into two parts. For the first part, we note:

$$\begin{aligned}
 2721 \quad \sum_{t=0}^{\infty} \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) \cdot 1 &= \sum_{t=0}^{\infty} \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) \cdot \pi_*(s_0, a_0, \dots, s, a) \\
 2722 &= \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0, \dots} \tau(s_0, a_0, \dots, s) \cdot \pi_*(s_0, a_0, \dots, s, a) \\
 2723 &= \frac{D^{\pi^*}(s, a)}{1 - \gamma}.
 \end{aligned} \tag{92}$$

2729 For the second part, we similarly compute:

$$\begin{aligned}
 2731 \quad \sum_{t=0}^{\infty} (t + 1) \gamma^t \sum_{\substack{s_0, a_0, \dots \\ \pi_*\text{-comp}}} \tau(s_0, a_0, \dots, s) &= \sum_{t=0}^{\infty} \frac{\partial}{\partial \gamma} \gamma^{t+1} P(s_t = s, a_t = a | \pi_*) \\
 2732 &= \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1 - \gamma} \cdot D^{\pi^*}(s, a) \right).
 \end{aligned} \tag{93}$$

2736 Putting Equations (92) and (93) into Equation (91) gives the first equation of Equation (89) for the
 2737 case that $(s, a) \in \text{supp } D^{\pi^*}$. For the case that $(s, a) \notin \text{supp } D^{\pi^*}(s, a)$, the inequality is trivial since
 2738 then $D^{\pi^*}(s, a) = 0$ and since the stated derivative is easily shown to be non-negative by writing out
 2739 the occupancy explicitly (i.e., by reversing the previous computation).

2740 This then implies

$$\begin{aligned}
 2741 \quad D^\pi(\text{supp } D^{\pi^*}) &= \sum_{(s, a) \in \text{supp } D^{\pi^*}} D^\pi(s, a) \\
 2742 &\geq \sum_{(s, a) \in \text{supp } D^{\pi^*}} \left(D^{\pi^*}(s, a) - \delta \cdot (1 - \gamma) \cdot \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1 - \gamma} D^{\pi^*}(s, a) \right) \right) \\
 2743 &= 1 - \delta \cdot (1 - \gamma) \cdot \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1 - \gamma} \sum_{(s, a) \in \text{supp } D^{\pi^*}} D^{\pi^*}(s, a) \right) \\
 2744 &= 1 - \delta \cdot (1 - \gamma) \cdot \frac{1}{(1 - \gamma)^2} \\
 2745 &= 1 - \frac{\delta}{1 - \gamma}.
 \end{aligned}$$

This shows the first inequality in Equation (90). To show the second inequality in Equation (89), we use the first one and compute:

$$\begin{aligned}
D^\pi(s, a) &= 1 - \sum_{(s', a') \neq (s, a)} D^\pi(s', a') \\
&\leq 1 - \sum_{(s', a') \in \text{supp } D^{\pi_*} \setminus \{(s, a)\}} D^\pi(s', a') \\
&\leq 1 - \sum_{(s', a') \in \text{supp } D^{\pi_*} \setminus \{(s, a)\}} D^{\pi_*}(s', a') \\
&\quad + \sum_{(s', a') \in \text{supp } D^{\pi_*} \setminus \{(s, a)\}} \delta \cdot (1 - \gamma) \cdot \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1 - \gamma} D^{\pi_*}(s', a') \right) \\
&\leq D^{\pi_*}(s, a) + \frac{\delta}{1 - \gamma},
\end{aligned}$$

where in the last step we again used the trick of the previous computation of pulling the sum through the derivative. Finally, we prove the second inequality in Equation (90), using what we know so far. First, note that

$$\delta \cdot (1 - \gamma) \cdot \frac{\partial}{\partial \gamma} \left(\frac{\gamma}{1 - \gamma} D^{\pi_*}(s, a) \right) \leq \frac{\delta}{1 - \gamma}$$

since we showed that the left-hand-side is non-negative and sums to the right-hand-side over all (s, a) . Consequently, we obtain:

$$\begin{aligned}
\|D^\pi - D^{\pi_*}\| &= \sqrt{\sum_{(s, a)} (D^\pi(s, a) - D^{\pi_*}(s, a))^2} \\
&\leq \sqrt{\sum_{(s, a)} \left| \frac{\delta}{1 - \gamma} \right|^2} \\
&= \sqrt{|\mathcal{S} \times \mathcal{A}|} \cdot \frac{\delta}{1 - \gamma}.
\end{aligned}$$

This finishes the proof. \square

We now fix more constants and notation. Define $\mathcal{S}_0 := \text{supp } \mu_0$ as the support of μ_0 , and more generally \mathcal{S}_t as the states reachable within t timesteps using the fixed worst-case policy π_* :

$$\mathcal{S}_t := \left\{ s \mid \exists \pi_*\text{-compatible sequence } s_0, a_0, \dots, s_{k-1}, a_{k-1}, s \text{ for } k \leq t \right\}.$$

Since there are only finitely many states and $\mathcal{S}_t \subseteq \mathcal{S}_{t+1}$, there is a t_0 such that \mathcal{S}_{t_0} is maximal. Set $D^{\pi_*}(s) := \sum_a D^{\pi_*}(s, a)$. Recall the notation τ from Lemma C.35. Define the following constant which, given the MDP, only depends on $\delta > 0$ and π_* :

$$C(\delta, \pi_*, \mu_0, \tau, \gamma) := \min_{\substack{t \in [0: t_0] \\ s_0, a_0, \dots, s_{t-1}, a_{t-1}, s: \pi_*\text{-comp}}} \gamma^t \tau(s_0, a_0, \dots, s) \cdot (1 - \delta)^t \cdot \delta > 0. \quad (94)$$

We get the following result:

Lemma C.37. Define the reward function $\hat{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as follows:

$$\hat{R}(s, a) := \begin{cases} R(s, a), & (s, a) \notin \text{supp } D^{\pi_*}, \\ \max R + \frac{\lambda}{C(\delta, \pi_*, \mu_0, \tau, \gamma)} \cdot \omega(\pi_*), & \text{else.} \end{cases} \quad (95)$$

Assume that $\hat{\pi}$ is (λ, ω) -RLHF optimal with respect to \hat{R} . Then for all $(s, a) \in \text{supp } D^{\pi_*}$, we have $\hat{\pi}(a \mid s) \geq 1 - \delta$.

Proof. We show this statement by induction over the number of timesteps that π_* needs to reach a given state. Thus, first assume $s \in \mathcal{S}_0$ and $a = \pi_*(s)$. We do a proof by contradiction. Thus, assume that $\hat{\pi}(a | s) < 1 - \delta$. This means that $\sum_{a' \neq a} \hat{\pi}(a' | s) \geq \delta$, and consequently

$$\sum_{a' \neq a} D^{\hat{\pi}}(s, a') \geq \mu_0(s) \cdot \delta \geq C(\delta, \pi_*, \mu_0, \tau, \gamma). \quad (96)$$

We now claim that from this it follows that π_* is more optimal than $\hat{\pi}$ with respect to RLHF, a contradiction to the optimality of $\hat{\pi}$. Indeed:

$$\begin{aligned} J_{\hat{R}}(\hat{\pi}) - \lambda \omega(\hat{\pi}) &\stackrel{(1)}{\leq} J_{\hat{R}}(\hat{\pi}) \\ &\stackrel{(2)}{=} \sum_{a' \neq a} D^{\hat{\pi}}(s, a') \cdot R(s, a') + \sum_{(s', a') \notin \{s\} \times \mathcal{A} \setminus \{a\}} D^{\hat{\pi}}(s', a') \cdot \hat{R}(s', a') \\ &\stackrel{(3)}{\leq} \sum_{a' \neq a} D^{\hat{\pi}}(s, a') \cdot \max R + \hat{R}(s, a) \cdot \sum_{(s', a') \notin \{s\} \times \mathcal{A} \setminus \{a\}} D^{\hat{\pi}}(s', a') \\ &= \sum_{a' \neq a} D^{\hat{\pi}}(s, a') \cdot \max R + \left(1 - \sum_{a' \neq a} D^{\hat{\pi}}(s, a')\right) \cdot \hat{R}(s, a) \\ &\stackrel{(4)}{\leq} C(\delta, \pi_*, \mu_0, \tau, \gamma) \cdot \max R + (1 - C(\delta, \pi_*, \mu_0, \tau, \gamma)) \cdot \hat{R}(s, a) \\ &\stackrel{(5)}{=} J_{\hat{R}}(\pi_*) + C(\delta, \pi_*, \mu_0, \tau, \gamma) \cdot (\max R - \hat{R}(s, a)) \\ &\stackrel{(6)}{=} J_{\hat{R}}(\pi_*) - C(\delta, \pi_*, \mu_0, \tau, \gamma) \cdot \frac{\lambda}{C(\delta, \pi_*, \mu_0, \tau, \gamma)} \cdot \omega(\pi_*) \\ &= J_{\hat{R}}(\pi_*) - \lambda \omega(\pi_*). \end{aligned} \quad (97)$$

In step (1), we use the non-negativity of ω . In step (2), we use that $(s, a') \notin \text{supp } D^{\pi_*}$, and so $\hat{R}(s, a') = R(s, a')$. In the right term in step (3), we use that $(s, a) \in \text{supp } D^{\pi_*}$, and thus $\hat{R}(s, a) \geq \hat{R}(s', a')$, by definition of \hat{R} . In step (4), we use that $\hat{R}(s, a) \geq \max R$ and Equation (96). Step (5) uses that $J_{\hat{R}}(\pi_*) = \hat{R}(s, a)$, following from the fact that \hat{R} is constant for policy π_* . Step (6) uses the concrete definition of \hat{R} . Thus, we have showed a contradiction to the RLHF-optimality of $\hat{\pi}$, from which it follows that $\hat{\pi}(a | s) \geq 1 - \delta$.

Now assume the statement is already proven for $t - 1$ and let $s \in \mathcal{S}_t \setminus \mathcal{S}_{t-1}$. Then there exists a π_* -compatible sequence $s_0, a_0, \dots, s_{t-1}, a_{t-1}$ leading to s . We necessarily have $s_i \in \mathcal{S}_i$ for all $i = 0, \dots, t-1$, and so we obtain $\hat{\pi}(a_i | s_i) \geq 1 - \delta$ by the induction hypothesis. Now, let $a := \pi_*(s)$ and assume we had $\hat{\pi}(a | s) < 1 - \delta$. As before, we then have $\sum_{a' \neq a} \hat{\pi}(a' | s) \geq \delta$. Consequently, we get

$$\begin{aligned} \sum_{a' \neq a} D^{\hat{\pi}}(s, a') &\geq \gamma^t \cdot \sum_{a' \neq a} \tau(s_0, a_0, \dots, s) \cdot \hat{\pi}(s_0, a_0, \dots, s, a') \\ &\geq \gamma^t \cdot \tau(s_0, a_0, \dots, s) \cdot (1 - \delta)^t \cdot \delta \\ &\geq C(\delta, \pi_*, \mu_0, \tau, \gamma) \end{aligned}$$

Then the same computation as in Equation (97) leads to the same contradiction again, and we are done. \square

Theorem C.38. *Define*

$$\delta := \frac{(1 - \gamma) \cdot (1 - L) \cdot \text{range } J_R}{\sqrt{|\mathcal{S} \times \mathcal{A}|} \cdot \|R\|} > 0.$$

Let $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be our MDP. Set

$$C := C(\mathcal{M}, \pi_*, L, \lambda, \omega) := \frac{\lambda \cdot \omega(\pi_*)}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)} < \infty, \quad (98)$$

with the “inner” $C(\delta, \pi_*, \mu_0, \tau, \gamma)$ defined in Equation (94). Assume that

$$D(\text{supp } D^{\pi_*}) \leq \frac{\epsilon}{1 + C}. \quad (99)$$

Then $D \in \text{unsafe}(R, \epsilon, L, \lambda, \omega)$.

Proof. We prove the theorem by showing that for every data distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$ that fulfills the conditions of Theorem C.38, there exists a reward function \hat{R} together with a (λ, ω) -RLHF optimal policy $\hat{\pi}$ with respect to \hat{R} such that

- $\mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s,a) - R(s,a)|}{\text{range } R} \right] \leq \epsilon,$
- $\text{Reg}^R(\hat{\pi}) \geq L.$

Towards that goal, define \hat{R} as in Equation (95) and $\hat{\pi}$ as a (λ, ω) -RLHF optimal policy for \hat{R} . Then Lemma C.37 shows that $\hat{\pi}(s | a) \geq 1 - \delta$ for all $(s, a) \in \text{supp } D^{\pi_*}$. Consequently, Lemma C.36 implies that

$$\|D^{\hat{\pi}} - D^{\pi_*}\| \leq \sqrt{|\mathcal{S} \times \mathcal{A}|} \cdot \frac{\delta}{1 - \gamma} = \frac{(1 - L) \cdot \text{range } J_R}{\|R\|}.$$

Consequently, Lemma C.34 shows that $\text{Reg}^R(\hat{\pi}) \geq L$, and thus the second claim. For the first claim, note that

$$\begin{aligned} \mathbb{E}_{(s,a) \sim D} \left[|\hat{R}(s,a) - R(s,a)| \right] &= \sum_{(s,a) \in \text{supp } D^{\pi_*}} D(s,a) \cdot \left(\max R + \frac{\lambda}{C(\delta, \pi_*, \mu_0, \tau, \gamma)} \omega(\pi_*) - R(s,a) \right) \\ &\leq D(\text{supp } D^{\pi_*}) \cdot \left(\text{range } R + \frac{\lambda}{C(\delta, \pi_*, \mu_0, \tau, \gamma)} \omega(\pi_*) \right) \\ &\leq \epsilon \cdot \text{range } R, \end{aligned}$$

where the last claim follows from the assumed inequality in $D(\text{supp } D^{\pi_*})$. \square

We obtain the following corollary, which is very similar to Proposition C.5. The main difference is that the earlier result only assumed a poly of regret L and not regret 1:

Corollary C.39. *Theorem C.38 specializes as follows for the case $\lambda = 0$: Assume $D(\text{supp } D^{\pi_*}) \leq \epsilon$. Then there exists a reward function \hat{R} together with an optimal policy $\hat{\pi}$ that satisfies the two inequalities from the previous result.*

Proof. This directly follows from $\lambda = 0$. For completeness, we note that the definition of \hat{R} also simplifies, namely to

$$\hat{R}(s, a) = \begin{cases} R(s, a), & (s, a) \notin \text{supp } D^{\pi_*} \\ \max R, & \text{else.} \end{cases}$$

\square

We now present another specialization of Theorem C.38. Namely, from now on, assume that $D = D^{\pi_{\text{ref}}}$ and $\omega(\pi) = \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}})$. In other words, the dataset used to evaluate the reward function is sampled from the same (safe) policy used in KL-regularization. This leads to the following condition specializing the one from Equation (99):

$$D^{\pi_{\text{ref}}}(\text{supp } D^{\pi_*}) \leq \frac{\epsilon}{1 + \frac{\lambda \cdot \mathbb{D}_{\text{KL}}(\pi_* || \pi_{\text{ref}})}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)}}. \quad (100)$$

π_{ref} now appears on both the left and right side of the equation, and so one can wonder whether it is ever possible that the inequality holds. After all, if $D^{\pi_{\text{ref}}}(\text{supp } D^{\pi_*})$ “gets smaller”, then $\mathbb{D}_{\text{KL}}(\pi_* || \pi_{\text{ref}})$ should usually get “larger”. However, halving each of the probabilities $D^{\pi_{\text{ref}}}(s, a)$ for $(s, a) \in \text{supp } D^{\pi_*}$ leads to only an increase by the addition of $\log 2$ of $\mathbb{D}_{\text{KL}}(\pi_* || \pi_{\text{ref}})$. Thus, intuitively, we expect the inequality to hold when the left-hand-side is very small. An issue is that the

KL divergence can disproportionately blow up in size if some *individual* probabilities $D^{\pi_{\text{ref}}}(s, a)$ for $(s, a) \in \text{supp } D^{\pi_*}$ are very small compared to other such probabilities. This can be avoided by a bound in the proportional difference of these probabilities. We thus obtain the following sufficient condition for a “negative result”:³

Corollary C.40. *Let the notation be as in Theorem C.38 and assume $D = D^{\pi_{\text{ref}}}$ and $\omega(\pi) = \mathbb{D}_{\text{KL}}(\pi || \pi_{\text{ref}})$. Let $K \geq 0$ be a constant such that*

$$\max_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a) \leq K \cdot \min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a).$$

Assume that

$$\min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a) \leq \left(\frac{\epsilon}{K \cdot |\mathcal{S}| \cdot \left(1 + \frac{\lambda}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)}\right)} \right)^2. \quad (101)$$

Then Equation (99) holds, and the conclusion of the theorem thus follows.

Proof. As argued before, the equation to show can be written as Equation (100). We can upper-bound the left-hand-side as follows:

$$\begin{aligned} D^{\pi_{\text{ref}}}(\text{supp } D^{\pi_*}) &= \sum_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a) \\ &\leq |\text{supp } D^{\pi_*}| \cdot \max_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a) \\ &\leq |\mathcal{S}| \cdot K \cdot \min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a). \end{aligned} \quad (102)$$

In one step, we used that π_* is assumed to be deterministic, which leads to a bound in the size of the support. Now, we lower-bound the other side by noting that

$$\begin{aligned} \mathbb{D}_{\text{KL}}(\pi_* || \pi_{\text{ref}}) &= \sum_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_*}(s, a) \cdot \log \frac{D^{\pi_*}(s, a)}{D^{\pi_{\text{ref}}}(s, a)} \\ &\leq \sum_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_*}(s, a) \cdot \log \frac{1}{\min_{(s',a') \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s', a')} \\ &= \log \frac{1}{\min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a)}. \end{aligned}$$

Thus, for the right-hand-side, we obtain

$$\frac{\epsilon}{1 + \frac{\lambda \cdot \mathbb{D}_{\text{KL}}(\pi_* || \pi_{\text{ref}})}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)}} \geq \frac{\epsilon}{1 + \frac{\lambda}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)} \cdot \log \frac{1}{\min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a)}} \quad (103)$$

Now, set $A := |\mathcal{S}| \cdot K$, $B := \frac{\lambda}{\text{range } R \cdot C(\delta, \pi_*, \mu_0, \tau, \gamma)}$ and $x := \min_{(s,a) \in \text{supp } D^{\pi_*}} D^{\pi_{\text{ref}}}(s, a)$. Then comparing with Equations (102) and (103), we are left with showing the following, which we also equivalently rewrite:

$$\begin{aligned} A \cdot x &\leq \frac{\epsilon}{1 + B \cdot \log \frac{1}{x}} \\ \iff A \cdot \left(x + Bx \log \frac{1}{x} \right) &\leq \epsilon. \end{aligned}$$

Now, together with the assumed condition on x from Equation (101), and upper-bounding the logarithm with a square-root, and x by \sqrt{x} since $x \leq 1$, we obtain:

$$\begin{aligned} A \cdot \left(x + Bx \log \frac{1}{x} \right) &\leq A \cdot (x + B\sqrt{x}) \\ &\leq A \cdot ((1 + B) \cdot \sqrt{x}) \\ &\leq A \cdot (1 + B) \cdot \frac{\epsilon}{A \cdot (1 + B)} \\ &= \epsilon. \end{aligned}$$

³The condition is quite strong and we would welcome attempts to weaken it.

2970 That was to show. □

2971

2972

2973

D REQUIREMENTS FOR SAFE OPTIMIZATION

2974

2975

2976

2977

In this section, we answer the question under which circumstances we can guarantee a safe optimization of a given reward function. Wherever applicable, we make the same assumptions as stated in Appendix C.1.

2978

2979

D.1 APPLYING BERGE’S MAXIMUM THEOREM

2980

2981

Definition D.1 (Correspondence). Let X, Y be two sets. A *correspondence* $C : X \rightrightarrows Y$ is a function $X \rightarrow \mathcal{P}(Y)$ from X to the power set of Y .

2982

2983

2984

Definition D.2 (Upper Hemicontinuous, Lower Hemicontinuous, Continuous, Compact-Valued). Let $C : X \rightrightarrows Y$ be a correspondence where X and Y are topological spaces. Then:

2985

2986

2987

- C is called *upper hemicontinuous* if for every $x \in X$ and every open set $V \subseteq Y$ with $C(x) \subseteq V$, there exists an open set $U \subseteq X$ with $x \in U$ and such that for all $x' \in U$ one has $C(x') \subseteq V$.

2988

2989

2990

2991

- C is called *lower hemicontinuous* if for every $x \in X$ and every open set $V \subseteq Y$ with $C(x) \cap V \neq \emptyset$, there exists an open set $U \subseteq X$ with $x \in U$ and such that for all $x' \in U$ one has $C(x') \cap V \neq \emptyset$.

2992

2993

2994

- C is called *continuous* if it is both upper and lower hemicontinuous.

2995

2996

2997

- C is called *compact-valued* if $C(x)$ is a compact subset of Y for all $x \in X$.

Theorem D.3 (Maximum Theorem, (Berge, 1963)). Let Θ and X be topological spaces, $f : \Theta \times X \rightarrow \mathbb{R}$ a continuous function, and $C : \Theta \rightrightarrows X$ be a continuous, compact-valued correspondence such that $C(\theta) \neq \emptyset$ for all $\theta \in \Theta$. Define the optimal value function $f^* : \Theta \rightarrow \mathbb{R}$ by

2998

2999

$$f^*(\theta) := \max_{x \in C(\theta)} f(\theta, x)$$

3000

and the maximizer function $C^* : \Theta \rightrightarrows X$ by

3001

3002

$$C^*(\theta) := \arg \max_{x \in C(\theta)} f(\theta, x) = \{x \in C(\theta) \mid f(\theta, x) = f^*(\theta)\}.$$

3003

3004

Then f^* is continuous and C^* is a compact-valued, upper hemicontinuous correspondence with nonempty values, i.e. $C^*(\theta) \neq \emptyset$ for all $\theta \in \Theta$.

3005

3006

3007

3008

3009

3010

3011

We now show that this theorem corresponds to our setting. Namely, replace X be by Π , the set of all policies. Every policy $\pi \in \Pi$ can be viewed as a vector $\vec{\pi} = (\pi(a \mid s))_{s \in \mathcal{S}, a \in \mathcal{A}} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$, and so we view Π as a subset of $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. Π inherits the standard Euclidean metric and thus topology from $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. Replace Θ by \mathcal{R} , the set of all reward functions. We can view each reward function $R \in \mathcal{R}$ as a vector $\vec{R} = (R(s, a))_{(s, a) \in \mathcal{S} \times \mathcal{A}} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$. So we view \mathcal{R} as a subset of $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ and thus a topological space. Replace f by the function $J : \mathcal{R} \times \Pi \rightarrow \mathbb{R}$ given by

3012

3013

$$J(R, \pi) := J^R(\pi) = \eta^\pi \cdot \vec{R}.$$

3014

3015

3016

Take as the correspondence $C : \mathcal{R} \rightrightarrows \Pi$ the trivial function $C(R) := \Pi$ that maps every reward function to the full set of policies.

3017

3018

3019

3020

Proposition D.4. These definitions satisfy the conditions of Theorem D.3, that is:

3021

3022

3023

1. $J : \mathcal{R} \times \Pi \rightarrow \mathbb{R}$ is continuous.

2. $C : \mathcal{R} \rightrightarrows \Pi$ is continuous and compact-valued with non-empty values.

Proof. Let us prove 1. Since the scalar product is continuous, it is enough to show that $\eta : \Pi \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ is continuous. Let $(s, a) \in \mathcal{S} \times \mathcal{A}$ be arbitrary. Then it is enough to show that each componentfunction $\eta(s, a) : \Pi \rightarrow \mathbb{R}$ given by

$$[\eta(s, a)](\pi) := \eta^\pi(s, a)$$

3024 is continuous.

3025 Now, for any $t \geq 0$, define the function $P_t(s, a) : \Pi \rightarrow \mathbb{R}$ by

$$3027 \quad [P_t(s, a)](\pi) := P(s_t = s, a_t = a \mid \xi \sim \pi).$$

3028 We obtain

$$3029 \quad \eta(s, a) = \sum_{t=0}^{\infty} \gamma^t P_t(s, a).$$

3032 Furthermore, this convergence is uniform since $[P_t(s, a)](\pi) \leq 1$ for all π and since $\sum_{t=0}^{\infty} \gamma^t$ is a
3033 convergent series. Thus, by the uniform limit theorem, it is enough to show that each $P_t(s, a)$ is a
3034 continuous function.

3035 Concretely, we have

$$3036 \quad [P_t(s, a)](\pi) = \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} P(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s, a \mid \xi \sim \pi)$$

$$3037 \quad = \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} \mu_0(s_0) \cdot \pi(a_0 \mid s_0) \cdot \left[\prod_{l=1}^{t-1} \tau(s_l \mid s_{l-1}, a_{l-1}) \cdot \pi(a_l \mid s_l) \right] \cdot \tau(s \mid s_{t-1}, a_{t-1}) \cdot \pi(a \mid s).$$

3042 Since \mathcal{S} and \mathcal{A} are finite, this whole expression can be considered as a polynomial with variables
3043 given by all $\pi(a \mid s)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and coefficients specified by μ_0 and τ . Since polynomials
3044 are continuous, this shows the result.

3045 Let us prove 2. Since $\Pi \neq \emptyset$, C has non-empty values. Furthermore, Π is compact because it is a
3046 finite cartesian product of compact simplices. And finally, since C is constant, it is easily seen to be
3047 continuous. That was to show. \square

3049 Define the optimal value function $J^* : \mathcal{R} \rightarrow \mathbb{R}$ by

$$3050 \quad J^*(R) := \max_{\pi \in \Pi} J^R(\pi)$$

3052 and the maximizer function $\Pi^* : \mathcal{R} \rightrightarrows \Pi$ by

$$3053 \quad \Pi^*(R) := \arg \max_{\pi \in \Pi} J^R(\pi) = \{\pi \in \Pi \mid J^R(\pi) = J^*(R)\}.$$

3055 **Corollary D.5.** *J^* is continuous and Π^* is upper hemicontinuous and compact-valued with non-*
3056 *empty values.*

3058 *Proof.* This follows from Theorem D.3 and Proposition D.4. \square

3060 In particular, every reward function has a compact and non-empty set of optimal policies, and their
3061 value changes continuously with the reward function. The most important part of the corollary is the
3062 upper hemicontinuity, which has the following consequence:

3063 **Corollary D.6.** *Let R be a fixed, non-trivial reward function, meaning that $\max J^R \neq \min J^R$. Let*
3064 *$U \in (0, 1]$ be arbitrary. Then there exists $\epsilon > 0$ such that for all $\hat{R} \in \mathcal{B}_\epsilon(R)$ and all $\hat{\pi} \in \Pi^*(\hat{R})$, we*
3065 *have $\text{Reg}^R(\hat{\pi}) < U$.*

3067 *Proof.* The condition $\max J^R \neq \min J^R$ ensures that the regret function $\text{Reg}^R : \Pi \rightarrow [0, 1]$ is
3068 well-defined. Recall its definition:

$$3070 \quad \text{Reg}^R(\pi) = \frac{\max J^R - J^R(\pi)}{\max J^R - \min J^R}.$$

3072 Since J^R is continuous by Proposition D.4, the regret function Reg^R is continuous as well. Conse-
3073 quently, the set $V := (\text{Reg}^R)^{-1}([0, U])$ is open in Π .

3074 Notice that $\Pi^*(R) \subseteq V$ (optimal policies have no regret). Thus, by Corollary D.5, there exists an
3075 open set $W \subseteq \mathcal{R}$ with $R \in W$ such that for all $\hat{R} \in W$ we have $\Pi^*(\hat{R}) \subseteq V$. Consequently, for
3076 all $\hat{\pi} \in \Pi^*(\hat{R})$, we get $\text{Reg}^R(\hat{\pi}) < U$. Since W is open, it contains a whole epsilon ball around R ,
3077 showing the result. \square

3078 Now we translate the results to the distance defined by D , a data distribution. Namely, let $D \in$
 3079 $\Delta(\mathcal{S} \times \mathcal{A})$ a distribution that assigns a positive probability to each transition. Then define the D -norm
 3080 by

$$3081 \quad d^D(R) := \mathbb{E}_{(s,a) \sim D} [|R(s,a)|].$$

3082 This is indeed a norm, i.e.: for all $\alpha \in \mathbb{R}$ and all $R, R' \in \mathcal{R}$, we have

- 3083 • $d^D(R + R') \leq d^D(R) + d^D(R')$;
- 3084 • $d^D(\alpha \cdot R) = |\alpha| \cdot d^D(R)$;
- 3085 • $d^D(R) = 0$ if and only if $R = 0$.

3086 For the third property, one needs the assumption that $D(s, a) > 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$.

3087 This norm then induces a metric that we denote the same way:

$$3088 \quad d^D(R, R') := d^D(R - R').$$

3089 We obtain:

3090 **Corollary D.7.** *Let $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ be an arbitrary non-trivial MDP, meaning that $\max J^R \neq$
 3091 $\min J^R$. Furthermore, let $L \in (0, 1]$ be arbitrary, and $D \in \Delta(\mathcal{S} \times \mathcal{A})$ a positive data distribution,
 3092 i.e., a distribution D such that $\forall (s, a) \in \mathcal{S} \times \mathcal{A}, D(s, a) > 0$. Then there exists $\epsilon > 0$ such that
 3093 $D \in \text{safe}(R, \epsilon, L)$*

3094 *Proof.* To prove the corollary, we will show that there exists $\epsilon > 0$ such that for all $\hat{R} \in \mathcal{R}$ with

$$3095 \quad \frac{d^D(R, \hat{R})}{\text{range } R} < \epsilon$$

3096 and all $\hat{\pi} \in \Pi^*(\hat{R})$ we have $\text{Reg}^R(\hat{\pi}) < L$. We know from Corollary D.6 that there is $\epsilon' > 0$ such
 3097 that for all $\hat{R} \in \mathcal{B}_{\epsilon'}(R)$ and all $\hat{\pi} \in \Pi^*(\hat{R})$, we have $\text{Reg}^R(\hat{\pi}) < L$. Now, let $c > 0$ be a constant
 3098 such that

$$3099 \quad c \cdot \|R' - R''\| \leq d^D(R', R'')$$

3100 for all $R', R'' \in \mathcal{R}$, where $\|\cdot\|$ is the standard Euclidean norm. This exists since all norms in $\mathbb{R}^{\mathcal{S} \times \mathcal{A}}$
 3101 are equivalent, but one can also directly argue that

$$3102 \quad c := \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a)$$

3103 is a valid choice. Then, set

$$3104 \quad \epsilon := \epsilon' \cdot \frac{c}{\text{range } R}.$$

3105 Then for all $\hat{R} \in \mathcal{R}$ with

$$3106 \quad \frac{d^D(R, \hat{R})}{\text{range } R} < \epsilon$$

3107 we obtain

$$3108 \quad \begin{aligned} \|R - \hat{R}\| &\leq \frac{d^D(R, \hat{R})}{c} \\ &= \frac{d^D(R, R')}{\text{range } R} \cdot \frac{\text{range } R}{c} \\ &\leq \epsilon \cdot \frac{\text{range } R}{c} \\ &= \epsilon'. \end{aligned}$$

3109 Thus, for all $\hat{\pi} \in \Pi^*(\hat{R})$, we obtain $\text{Reg}^R(\hat{\pi}) < L$, showing the result. \square

3110 **Remark D.8.** If $c := \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s, a)$ is very small, then the proof of the preceding corollary
 3111 shows that $d^D(R, \hat{R})$ must be correspondingly smaller to guarantee a low regret of $\hat{\pi} \in \Pi^*(\hat{R})$. This
 3112 makes sense since a large effective distance between R and \hat{R} can “hide” in the regions where D is
 3113 small when distance is measured via d^D .

3132 D.2 ELEMENTARY PROOF OF A REGRET BOUND
3133

3134 In this section, we provide another elementary proof of a regret bound, but without reference to
3135 Berge’s theorem. This will also lead to a better quantification of the bound. In an example, we will
3136 show that the bound we obtain is tight.

3137 Define the cosine of an angle between two vectors ad hoc as usual:

$$3138 \cos(\text{ang}(v, w)) := \frac{v \cdot w}{\|v\| \cdot \|w\|},$$

3139 where $v \cdot w$ is the dot product.

3140 **Lemma D.9.** *Let R, \hat{R} be two reward functions. Then for any policy π , we have*

$$3141 J^R(\pi) - J^{\hat{R}}(\pi) = \frac{1}{1-\gamma} \cdot \|D^\pi\| \cdot \|R - \hat{R}\| \cdot \cos(\text{ang}(\eta^\pi, \vec{R} - \vec{\hat{R}})).$$

3142 *Proof.* We have

$$3143 J^R(\pi) - J^{\hat{R}}(\pi) = \eta^\pi \cdot (\vec{R} - \vec{\hat{R}}) = \|\eta^\pi\| \cdot \|\vec{R} - \vec{\hat{R}}\| \cdot \cos(\text{ang}(\eta^\pi, \vec{R} - \vec{\hat{R}})).$$

3144 The result follows from $\eta^\pi = \frac{1}{1-\gamma} \cdot D^\pi$. □

3145 we will make use of another lemma:

3146 **Lemma D.10.** *Let a, \hat{a} , and r be three vectors. Assume $a \cdot \hat{a} \geq 0$, where \cdot is the dot product. Then*

$$3147 \cos(\text{ang}(a, r)) - \cos(\text{ang}(\hat{a}, r)) \leq \sqrt{2}.$$

3148 *Proof.* None of the angles change by replacing any of the vectors with a normed version. We can
3149 thus assume $\|a\| = \|\hat{a}\| = \|r\| = 1$. We obtain

$$3150 \begin{aligned} |\cos(\text{ang}(a, r)) - \cos(\text{ang}(\hat{a}, r))|^2 &= |a \cdot r - \hat{a} \cdot r|^2 \\ 3151 &= |(a - \hat{a}) \cdot r|^2 \\ 3152 &\leq \|a - \hat{a}\|^2 \cdot \|r\|^2 \\ 3153 &= \|a - \hat{a}\|^2 \\ 3154 &= \|a\|^2 + \|\hat{a}\|^2 - 2a \cdot \hat{a} \\ 3155 &\leq 2. \end{aligned}$$

3156 In the first, fourth, and sixth step, we used that all vectors are normed. In the third step, we used the
3157 Cauchy-Schwarz inequality. Finally, we used that $a \cdot \hat{a} \geq 0$. The result follows. □

3158 Recall that for two vectors v, w , the projection of v onto w is defined by

$$3159 \text{proj}_w v := \frac{v \cdot w}{\|w\|^2} w.$$

3160 This projection is a multiple of w , and it minimizes the distance to v :

$$3161 \|v - \text{proj}_w v\| = \min_{\alpha \in \mathbb{R}} \|v - \alpha w\|.$$

3162 We can now formulate and prove our main regret bound:

3163 **Theorem D.11.** *Let R be a fixed, non-trivial reward function, meaning that $\max J^R \neq \min J^R$.
3164 Then for all $\hat{R} \in \mathcal{R}$ and all $\hat{\pi} \in \Pi^*(\hat{R})$, we have*

$$3165 \text{Reg}^R(\hat{\pi}) \leq \frac{\sqrt{2}}{(1-\gamma) \cdot (\max J^R - \min J^R)} \cdot \|\vec{R} - \vec{\hat{R}}\|.$$

3166 Furthermore, if $\vec{R} \cdot \vec{\hat{R}} \geq 0$, then we also obtain the following stronger bound:

$$3167 \text{Reg}^R(\hat{\pi}) \leq \frac{\sqrt{2}}{(1-\gamma) \cdot (\max J^R - \min J^R)} \cdot \left\| \vec{R} - \text{proj}_{\vec{\hat{R}}} \vec{R} \right\|.$$

3186 Now, let $D \in \Delta(\mathcal{S} \times \mathcal{A})$ be a data distribution. Then we obtain the following consequence:
3187

$$3188 \text{Reg}^R(\hat{\pi}) \leq \frac{\sqrt{2}}{(1-\gamma) \cdot (\max J^R - \min J^R) \cdot \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a)} \cdot d^D(R, \hat{R}).$$

3190 *Proof.* We start with the first claim. First, notice that the inequality we want to show is equivalent to
3191 the following:
3192

$$3193 J^R(\hat{\pi}) \geq \max J^R - \frac{\sqrt{2}}{1-\gamma} \cdot \|\vec{R} - \vec{\hat{R}}\|. \quad (104)$$

3194 From Lemma D.9, we obtain

$$3196 J^R(\hat{\pi}) = J^{\hat{R}}(\hat{\pi}) + \frac{1}{1-\gamma} \cdot \|D^{\hat{\pi}}\| \cdot \|\vec{R} - \vec{\hat{R}}\| \cdot \cos(\text{ang}(\eta^{\hat{\pi}}, R - \hat{R})).$$

3198 Now, let $\pi \in \Pi^*(R)$ be an optimal policy for R . Then also from Lemma D.9, we obtain

$$3200 \max J^R = J^R(\pi) = J^{\hat{R}}(\pi) + \frac{1}{1-\gamma} \cdot \|D^\pi\| \cdot \|\vec{R} - \vec{\hat{R}}\| \cdot \cos(\text{ang}(\eta^\pi, R - \hat{R}))$$

$$3202 \leq J^{\hat{R}}(\hat{\pi}) + \frac{1}{1-\gamma} \cdot \|D^\pi\| \cdot \|\vec{R} - \vec{\hat{R}}\| \cdot \cos(\text{ang}(\eta^\pi, R - \hat{R})).$$

3204 In the last step, we used that $\hat{\pi} \in \Pi^*(\vec{R})$ and so $J^{\hat{R}}(\pi) \leq J^{\hat{R}}(\hat{\pi})$. Combining both computations, we
3205 obtain:

$$3206 J^R(\hat{\pi}) \geq \max J^R - \frac{1}{1-\gamma} \cdot \|\vec{R} - \vec{\hat{R}}\| \cdot \left[\|D^\pi\| \cdot \cos(\text{ang}(\eta^\pi, R - \hat{R})) - \|D^{\hat{\pi}}\| \cdot \cos(\text{ang}(\eta^{\hat{\pi}}, R - \hat{R})) \right]$$

3208 Since we want to show Equation (104), we are done if we can bound the big bracket by $\sqrt{2}$. By the
3209 Cauchy-Schwarz inequality, $\cos(\text{ang}(v, w)) \in [-1, 1]$ for all vectors v, w . Thus, if the first cosine
3210 term is negative or the second cosine term is positive, then since $\|D^\pi\| \leq \|D^\pi\|_1 = 1$, the bound by
3211 $\sqrt{2}$ is trivial. Thus, assume that the first cosine term is positive and the second is negative. We obtain

$$3213 \|D^\pi\| \cdot \cos(\text{ang}(\eta^\pi, R - \hat{R})) - \|D^{\hat{\pi}}\| \cdot \cos(\text{ang}(\eta^{\hat{\pi}}, R - \hat{R}))$$

$$3214 \leq \cos(\text{ang}(\eta^\pi, R - \hat{R})) - \cos(\text{ang}(\eta^{\hat{\pi}}, R - \hat{R}))$$

$$3215 \leq \sqrt{2}$$

3218 by Lemma D.10. Here, we used that η^π and $\eta^{\hat{\pi}}$ have only non-negative entries and thus also
3219 nonnegative dot product $\eta^\pi \cdot \eta^{\hat{\pi}} \geq 0$.

3220 For the second claim, notice the following: if $\vec{R} \cdot \vec{\hat{R}} \geq 0$, then $\text{proj}_{\vec{R}} \vec{\hat{R}} = \alpha \cdot \vec{R}$ for some constant
3221 $\alpha \geq 0$. Consequently, we have $\hat{\pi} \in \Pi^*(\text{proj}_{\vec{R}} \vec{\hat{R}})$. The claim thus follows from the first result.
3222

3224 For the third claim, notice that

$$3225 \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \cdot \|\vec{R} - \vec{\hat{R}}\| \leq \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \cdot \|\vec{R} - \vec{\hat{R}}\|_1$$

$$3226 = \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \cdot \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} |R(s,a) - \hat{R}(s,a)|$$

$$3227 \leq \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a) \cdot |R(s,a) - \hat{R}(s,a)|$$

$$3228 = d^D(R, \hat{R}).$$

3230 So the first result implies the third. \square

3235 *Remark D.12.* As one can easily see geometrically, but also prove directly, there is the following
3236 equality of sets for a reward function R

$$3237 \left\{ \text{proj}_{\vec{R}} \vec{R} \mid \hat{R} \in \mathcal{R} \right\} = \left\{ \frac{1}{2} \vec{R} + \frac{1}{2} \|\vec{R}\| v \mid v \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}, \|v\| = 1 \right\}.$$

3238 In other words, the projections form a sphere of radius $\frac{1}{2} \|\vec{R}\|$ around the midpoint $\frac{1}{2} \vec{R}$.
3239

3240 We now show that the regret bound is tight:

3241 **Example D.13.** Let $U \in [0, 1]$ and $\gamma \in [0, 1)$ be arbitrary. Then there exists an MDP
3242 $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$ together with a reward function \hat{R} with $\vec{R} \cdot \vec{\hat{R}} \geq 0$ and a policy $\hat{\pi} \in \Pi^*(\hat{R})$
3243 such that

$$3244 \quad U = \text{Reg}^R(\hat{\pi}) = \frac{\sqrt{2}}{(1-\gamma) \cdot (\max J^R - \min J^R)} \cdot \|\vec{R} - \text{proj}_{\vec{\hat{R}}} \vec{R}\|.$$

3245 Furthermore, there exists a data distribution $D \in \Delta(\mathcal{S} \times \mathcal{A})$ such that

$$3246 \quad \text{Reg}^R(\hat{\pi}) = \frac{1}{(1-\gamma) \cdot (\max J^R - \min J^R) \cdot \min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a)} \cdot d^D(R, \hat{R}).$$

3247 *Proof.* If $U = 0$ then $\hat{R} = R$ always works. If $U > 0$, then set $\mathcal{S} = \{\star\}$ and $\mathcal{A} = \{a, b, c\}$. This
3248 determines τ and μ_0 . Define $R(x) := R(\star, x, \star)$ for any action $x \in \mathcal{A}$. Let $R(a) > R(b)$ be arbitrary
3249 and set

$$3250 \quad R(c) := R(a) - \frac{R(a) - R(b)}{U} \leq R(b).$$

3251 Define

$$3252 \quad \hat{R}(a) := \hat{R}(b) := \frac{R(a) + R(b)}{2}, \quad \hat{R}(c) := R(c).$$

3253 For a policy π , define $\pi(x) := \pi(x \mid \star)$ for any action $x \in \mathcal{A}$ and set the policy $\hat{\pi}$ by $\hat{\pi}(b) = 1$.

3254 We obtain:

$$\begin{aligned} 3255 \quad \|\vec{R} - \vec{\hat{R}}\| &= \sqrt{(R(a) - \hat{R}(a))^2 + (R(b) - \hat{R}(b))^2 + (R(c) - \hat{R}(c))^2} \\ 3256 &= \frac{1}{2} \cdot \sqrt{(R(a) - R(b))^2 + (R(b) - R(a))^2} \\ 3257 &= \frac{1}{\sqrt{2}} \cdot (R(a) - R(b)) \\ 3258 &= U \cdot \frac{R(a) - R(c)}{\sqrt{2}} \\ 3259 &= U \cdot \frac{\max R - \min R}{\sqrt{2}} \\ 3260 &= U \cdot \frac{(1-\gamma) \cdot (\max J^R - \min J^R)}{\sqrt{2}}. \end{aligned}$$

3261 Furthermore, we have

$$\begin{aligned} 3262 \quad \text{Reg}^R(\hat{\pi}) &= \frac{\frac{1}{1-\gamma} \cdot R(a) - \frac{1}{1-\gamma} \cdot R(b)}{\frac{1}{1-\gamma} \cdot R(a) - \frac{1}{1-\gamma} \cdot R(c)} \\ 3263 &= U. \end{aligned}$$

3264 This shows

$$3265 \quad U = \text{Reg}^R(\hat{\pi}) = \frac{\sqrt{2}}{(1-\gamma) \cdot (\max J^R - \min J^R)} \cdot \|\vec{R} - \vec{\hat{R}}\|.$$

3266 We are done if we can show that $\text{proj}_{\vec{\hat{R}}} \vec{R} = \vec{\hat{R}}$. This is equivalent to

$$3267 \quad \vec{\hat{R}} \cdot \vec{R} = \|\vec{\hat{R}}\|^2,$$

3268 which is in turn equivalent to

$$3269 \quad \vec{\hat{R}} \cdot [\vec{R} - \vec{\hat{R}}] = 0.$$

3270 This can easily be verified.

3271 Finally, for the claim about the data distribution, simply set $D(a) = D(b) = D(c) = \frac{1}{3}$. Then one
3272 can easily show that

$$3273 \quad \sqrt{2} \cdot \|\vec{R} - \vec{\hat{R}}\| = R(a) - R(b) = \frac{d^D(R, \hat{R})}{\min_{(s,a) \in \mathcal{S} \times \mathcal{A}} D(s,a)}.$$

3274 That shows the result. \square

3294 D.3 SAFE OPTIMIZATION VIA APPROXIMATED CHOICE PROBABILITIES
 3295

3296 In this section, we will show that for any chosen upper regret bound U , there is an $\epsilon > 0$ s.t. if the
 3297 choice probabilities of \hat{R} are ϵ -close to those of R , the regret of an optimal policy for \hat{R} is bounded
 3298 by U .

3299 Assume a finite time horizon T . Trajectories are then given by $\xi = s_0, a_0, s_1, \dots, a_{T-1}, s_T$. Let Ξ
 3300 be the set of all trajectories of length T . Let $D \in \Delta(\Xi)$ be a distribution. Assume that the human has
 3301 a true reward function R and makes choices in trajectory comparisons given by
 3302

$$3303 P_R(1 \mid \xi_1, \xi_2) = \frac{\exp(G(\xi_1))}{\exp(G(\xi_1)) + \exp(G(\xi_2))}. \quad (105)$$

3305 Here, the return function G is given by

$$3306 G(\xi) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t, s_{t+1}).$$

3309 We can then define the choice distance of proxy reward \hat{R} to true reward R as

$$3311 d_{\text{KL}}^D(R, \hat{R}) := \mathbb{E}_{\xi_1, \xi_2 \sim D \times D} \left[D_{\text{KL}} \left(P_R(\cdot \mid \xi_1, \xi_2) \parallel P_{\hat{R}}(\cdot \mid \xi_1, \xi_2) \right) \right]$$

3313 Here, $D_{\text{KL}} \left(P_R(\cdot \mid \xi_1, \xi_2) \parallel P_{\hat{R}}(\cdot \mid \xi_1, \xi_2) \right)$ is the Kullback-Leibler divergence of two binary
 3314 distributions over values 1, 2. Explicitly, for $P := P_R(\cdot \mid \xi_1, \xi_2)$ and similarly \hat{P} , we have
 3315

$$3316 D_{\text{KL}}(P \parallel \hat{P}) = P(1) \log \frac{P(1)}{\hat{P}(1)} + (1 - P(1)) \log \frac{1 - P(1)}{1 - \hat{P}(1)} \quad (106)$$

$$3318 = - \left[P(1) \log \hat{P}(1) + (1 - P(1)) \log (1 - \hat{P}(1)) \right] - H(P(1)).$$

3320 Here, $H(p) := -[p \log p + (1 - p) \log(1 - p)]$ is the binary entropy function.

3322 Fix in this whole section the true reward function R with $\max J^R \neq \min J^R$ in a fixed MDP.

3323 The goal of this section is to prove the following proposition:

3324 **Proposition D.14.** *Let $U \in (0, 1]$. Then there exists an $\epsilon > 0$ such that for all \hat{R} with*

$$3326 d_{\text{KL}}^D(R, \hat{R}) < \epsilon$$

3328 *and all $\hat{\pi} \in \Pi^*(\hat{R})$ we have $\text{Reg}^R(\hat{\pi}) < U$.*

3329 We prove this by chaining together four lemmas. The first of the four lemmas needs its own lemma,
 3330 so we end up with five lemmas overall:

3331 **Lemma D.15.** *Assume R, \hat{R} are two reward functions and π a policy. Then*

$$3333 |J^R(\pi) - J^{\hat{R}}(\pi)| \leq \max_{\xi \in \Xi} |G(\xi) - \hat{G}(\xi)|.$$

3335 *Proof.* We have

$$3336 |J^R(\pi) - J^{\hat{R}}(\pi)| = |\tilde{D}^\pi \cdot (G - \hat{G})|$$

$$3338 = \left| \sum_{\xi \in \Xi} \tilde{D}^\pi(\xi) \cdot (G(\xi) - \hat{G}(\xi)) \right|$$

$$3339 \leq \sum_{\xi \in \Xi} \tilde{D}^\pi(\xi) \cdot |G(\xi) - \hat{G}(\xi)|$$

$$3341 \leq \max_{\xi \in \Xi} |G(\xi) - \hat{G}(\xi)| \cdot \sum_{\xi \in \Xi} \tilde{D}^\pi(\xi)$$

$$3342 = \max_{\xi \in \Xi} |G(\xi) - \hat{G}(\xi)|$$

3347 In the last step, we used that distributions sum to one. □

3348 **Lemma D.16.** Let $U \in (0, 1]$. Then there exists $\sigma(U) > 0$ such that for all \hat{R} and $\hat{\pi} \in \Pi^*(\hat{R})$ for
 3349 which there exists $c \in \mathbb{R}$ such that $\max_{\xi \in \Xi} |\hat{G}(\xi) - G(\xi) - c| < \sigma(U)$, we have $\text{Reg}^R(\hat{\pi}) < U$.
 3350

3351 Concretely, we can set $\sigma(U) := \frac{\max J^R - \min J^R}{2} \cdot U$.
 3352

3353 *Proof.* Set $\sigma(U)$ as stated and let \hat{R} , $\hat{\pi}$ and c have the stated properties. The regret bound we want to
 3354 show is equivalent to the following statement:
 3355

$$3356 J^R(\hat{\pi}) > \max J^R - (\max J^R - \min J^R) \cdot U = \max J^R - 2\sigma(U). \quad (107)$$

3357 Let \tilde{c} be the constant such that $\hat{G} - c$ is the return function of $\hat{R} - \tilde{c}$. Concretely, one can set
 3358 $\tilde{c} = \frac{1-\gamma}{1-\gamma^T+1} \cdot c$. Lemma D.15 ensures that
 3359

$$3360 J^R(\hat{\pi}) > J^{\hat{R}-\tilde{c}}(\hat{\pi}) - \sigma(U). \quad (108)$$

3361 Now, let π be an optimal policy for R . Again, Lemma D.15 ensures
 3362

$$3363 \max J^R = J^R(\pi) < J^{\hat{R}-\tilde{c}}(\pi) + \sigma(U) \leq J^{\hat{R}-\tilde{c}}(\hat{\pi}) + \sigma(U). \quad (109)$$

3364 In the last step, we used that $\hat{\pi}$ is optimal for \hat{R} and thus also $\hat{R} - \tilde{c}$. Combining Equations (108)
 3365 and (109), we obtain the result, Equation (107). \square
 3366

3367 **Lemma D.17.** For $q \in (0, 1)$, define $g_q : (-q, 1 - q) \rightarrow \mathbb{R}$ by
 3368

$$3369 g_q(x) := \log \frac{q+x}{1-(q+x)}.$$

3370 Then for all $\sigma > 0$ there exists $\delta(q, \sigma) > 0$ such that for all $x \in (-q, 1 - q)$ with $|x| < \delta(q, \sigma)$, we
 3371 have $|g_q(x) - g_q(0)| < \sigma$.
 3372

3373 Concretely, one can choose
 3374

$$3375 \delta(q, \sigma) := (\exp(\sigma) - 1) \cdot \min \left\{ \frac{1}{\frac{1}{q} + \frac{\exp(\sigma)}{1-q}}, \frac{1}{\frac{1}{1-q} + \frac{\exp(\sigma)}{q}} \right\}$$

3377 *Proof.* If one does not care about the precise quantification, then the result is simply a reformulation
 3378 of the continuity of g_q at the point $x_0 = 0$.
 3379

3380 Now we show more specifically that $\delta(q, \sigma)$, as defined above, has the desired property. Namely,
 3381 notice the following sequence of equivalences (followed by a one-sided implication) that holds
 3382 whenever $x \geq 0$:

$$3383 |g_q(x) - g_q(0)| < \sigma \iff \log \frac{(q+x) \cdot (1-q)}{(1-(q+x)) \cdot q} < \sigma$$

$$3384 \iff \frac{(q+x) \cdot (1-q)}{(1-(q+x)) \cdot q} < \exp(\sigma)$$

$$3385 \iff (q+x) < (1-q-x) \cdot \frac{q}{1-q} \cdot \exp(\sigma)$$

$$3386 \iff \left(1 + \frac{q}{1-q} \cdot \exp(\sigma)\right) \cdot x < q \cdot (\exp(\sigma) - 1)$$

$$3387 \iff x < \frac{\exp(\sigma) - 1}{\frac{1}{q} + \frac{\exp(\sigma)}{1-q}}$$

$$3388 \iff |x| < \delta(q, \sigma).$$

3396 In the first step, we used the monotonicity of g_q to get rid of the absolute value. Similarly, whenever
 3397 $x \leq 0$, we have
 3398

$$3399 |g_q(x) - g_q(0)| < \sigma \iff x > \frac{1 - \exp(\sigma)}{\frac{1}{1-q} + \frac{\exp(\sigma)}{q}}$$

$$3400 \iff |x| < \delta(q, \sigma).$$

3401 This shows the result. \square

3402 **Lemma D.18.** For $q \in (0, 1)$, define $f_q : (0, 1) \rightarrow \mathbb{R}$ by

$$3403 \quad f_q(p) := -[q \log p + (1 - q) \log(1 - p)].$$

3404 Then for all $\delta > 0$ there exists $\mu(\delta) > 0$ such that for all $p \in (0, 1)$ with $f_q(p) < H(q) + \mu(\delta)$, we
3405 have $|p - q| < \delta$. Concretely, one can choose $\mu(\delta) := 2\delta^2$.

3406 *Proof.* Let $\delta > 0$ and define $\mu(\delta) := 2\delta^2$. Assume that $f_q(p) < H(q) + \mu(\delta)$. By Pinker's inequality,
3407 we have

$$3408 \quad \begin{aligned} 2(p - q)^2 &\leq q \log \frac{q}{p} + (1 - q) \cdot \log \frac{1 - q}{1 - p} \\ 3409 &= -H(q) + f_q(p) \\ 3410 &< \mu(\delta) \\ 3411 &= 2\delta^2. \end{aligned}$$

3412 Consequently, we have $|p - q| < \delta$. □

3413 **Lemma D.19.** Define $f_q(p)$ as in Lemma D.18. Then for all $\mu > 0$ there exists $\epsilon(\mu) > 0$ such that
3414 for all \hat{R} with $d_{\text{KL}}^D(R, \hat{R}) < \epsilon(\mu)$, we have the following for all $\xi_1, \xi_2 \in \Xi$:

$$3415 \quad f_{P_{R(1|\xi_1, \xi_2)}}(P_{\hat{R}}(1 | \xi_1, \xi_2)) < H(P_{R(1 | \xi_1, \xi_2)}) + \mu.$$

3416 Concretely, we can set $\epsilon(\mu) := \mu \cdot \min_{\xi_1, \xi_2 \in \Xi} D(\xi_1) \cdot D(\xi_2)$

3417 *Proof.* We have the following for all $\xi_1, \xi_2 \in \Xi$:

$$3418 \quad \begin{aligned} \mu \cdot \min_{\xi, \xi'} D(\xi) \cdot D(\xi') &= \epsilon(\mu) \\ 3419 &> d_{\text{KL}}^D(R, \hat{R}) \\ 3420 &= \mathbb{E}_{\xi, \xi' \sim D \times D} \left[D_{\text{KL}} \left(P_R(\cdot | \xi, \xi') \parallel P_{\hat{R}}(\cdot | \xi, \xi') \right) \right] \\ 3421 &\geq \left(\min_{\xi, \xi'} D(\xi) \cdot D(\xi') \right) \cdot D_{\text{KL}} \left(P_R(\cdot | \xi_1, \xi_2) \parallel P_{\hat{R}}(\cdot | \xi_1, \xi_2) \right) \end{aligned}$$

3422 Now, Equation (106) shows that

$$3423 \quad D_{\text{KL}} \left(P_R(\cdot | \xi_1, \xi_2) \parallel P_{\hat{R}}(\cdot | \xi_1, \xi_2) \right) = f_{P_{R(1|\xi_1, \xi_2)}}(P_{\hat{R}}(1 | \xi_1, \xi_2)) - H(P_{R(1 | \xi_1, \xi_2)}).$$

3424 The result follows. □

3425 **Corollary D.20.** Let $\sigma > 0$. Then there exists $\epsilon := \epsilon(\sigma) > 0$ such that $d_{\text{KL}}^D(R, \hat{R}) < \epsilon$ implies that
3426 there exists $c \in \mathbb{R}$ such that $\|G - (\hat{G} - c)\|_\infty < \sigma$.

3427 *Proof.* Set

$$3428 \quad \delta := \min_{\xi_1, \xi_2 \in \Xi \times \Xi} \delta \left(P_R(1 | \xi_1, \xi_2), \sigma \right), \quad \mu := \mu(\delta), \quad \epsilon := \epsilon(\mu),$$

3429 with the constants satisfying the properties from Lemmas D.17, D.18, and D.19. Now, let \hat{R} be such
3430 that $d_{\text{KL}}^D(R, \hat{R}) < \epsilon$.

3431 First of all, Lemma D.19 ensures that

$$3432 \quad f_{P_{R(1|\xi_1, \xi_2)}}(P_{\hat{R}}(1 | \xi_1, \xi_2)) < H(P_{R(1 | \xi_1, \xi_2)}) + \mu$$

3433 for all $\xi_1, \xi_2 \in \Xi$. Then Lemma D.18 shows that

$$3434 \quad |P_{\hat{R}}(1 | \xi_1, \xi_2) - P_R(1 | \xi_1, \xi_2)| < \delta$$

3435 for all $\xi_1, \xi_2 \in \Xi$. From Lemma D.17, we obtain that

$$3436 \quad \left| g_{P_{R(1|\xi_1, \xi_2)}} \left(P_{\hat{R}}(1 | \xi_1, \xi_2) - P_R(1 | \xi_1, \xi_2) \right) - g_{P_{R(1|\xi_1, \xi_2)}}(0) \right| < \sigma \quad (110)$$

3456 for all $\xi_1, \xi_2 \in \Xi$. Now, note that

$$3457 \quad g_{P_R(1|\xi_1, \xi_2)} \left(P_{\hat{R}}(1 | \xi_1, \xi_2) - P_R(1 | \xi_1, \xi_2) \right) = g_{P_{\hat{R}}(1|\xi_1, \xi_2)}(0).$$

3460 Furthermore, for $R' \in \{R, \hat{R}\}$, Equation (105) leads to the following computation:

$$3461 \quad g_{P_{R'}(1|\xi_1, \xi_2)}(0) = \log \frac{P_{R'}(1 | \xi_1, \xi_2)}{P_{R'}(2 | \xi_1, \xi_2)} \\ 3462 \quad = \log \frac{\exp(G'(\xi_1))}{\exp(G'(\xi_2))} \\ 3463 \quad = G'(\xi_1) - G'(\xi_2).$$

3468 Therefore, Equation (110) results in

$$3469 \quad \left| (\hat{G}(\xi_1) - G(\xi_1)) - (\hat{G}(\xi_2) - G(\xi_2)) \right| = \left| (\hat{G}(\xi_1) - \hat{G}(\xi_2)) - (G(\xi_1) - G(\xi_2)) \right| < \sigma$$

3470 for all $\xi_1, \xi_2 \in \Xi$. Now, let $\xi^* \in \Xi$ be any reference trajectory. Define $c := \hat{G}(\xi^*) - G(\xi^*)$. Then

3473 the preceding equation shows that

$$3474 \quad \left| \hat{G}(\xi) - G(\xi) - c \right| < \sigma$$

3475 for all $\xi \in \Xi$. That shows the claim. \square

3476 *Proof of Proposition D.14.* We prove Proposition D.14 by chaining together the constants from the
3477 preceding results. We have $U \in (0, 1]$ given. Then, set $\sigma := \sigma(U)$ and $\epsilon := \epsilon(\sigma)$ as in Lemma D.16
3478 and Corollary D.20. Now, let \hat{R} be such that $d_{\text{KL}}^D(R, \hat{R}) < \epsilon$ and let $\hat{\pi} \in \Pi^*(\hat{R})$. Our goal is to show
3479 that $\text{Reg}^R(\hat{\pi}) < U$.

3480 By Corollary D.20, there is $c > 0$ such that $\max_{\xi \in \Xi} |\hat{G}(\xi) - G(\xi) - c| < \sigma$. Consequently,
3481 Lemma D.16 ensures that $\text{Reg}^R(\hat{\pi}) < U$. This was to show. \square

3485 D.4 POSITIVE RESULT FOR REGULARIZED RLHF

3486 Here, we present simple positive results for regularized RLHF, both in a version with the expected
3487 reward distance, and in a version using the distance in choice probabilities. Some of it will directly
3488 draw from the positive results proved before.

3489 **Theorem D.21.** *Let $\lambda \in (0, \infty)$ be given and fixed. Assume we are given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$,
3490 and a data distribution $D \in \mathcal{S} \times \mathcal{A}$ which assigns positive probability to all transitions, i.e., $\forall (s, a) \in$
3491 $\mathcal{S} \times \mathcal{A}$, $D(s, a) > 0$. Let $\omega : \Pi \rightarrow \mathbb{R}$ be a continuous regularization function that has a reference
3492 policy π_{ref} as one of its minima.⁴ Assume that π_{ref} is not (λ, ω) -optimal for R and let $L =$
3493 $\text{Reg}^R(\pi_{\text{ref}})$. Then there exists $\epsilon > 0$ such that $D \in \text{safe}(R, \epsilon, L, \lambda, \omega)$.*

3494 *Proof.* We prove the theorem by showing that for every $D \in \Delta(\mathcal{S} \times \mathcal{A})$ such that $D(s, a) > 0$ for
3495 all $(s, a) \in \mathcal{S} \times \mathcal{A}$, there exists $\epsilon > 0$ such that for all \hat{R} with $\mathbb{E}_{(s, a) \sim D} \left[\frac{|\hat{R}(s, a) - R(s, a)|}{\text{range } R} \right] < \epsilon$ and
3496 all policies $\hat{\pi}$ that are (λ, ω) -RLHF optimal wrt. \hat{R} , we have $\text{Reg}^R(\hat{\pi}) < \text{Reg}^R(\pi_{\text{ref}})$. Because
3497 $L = \text{Reg}^R(\hat{\pi}) < \text{Reg}^R(\pi_{\text{ref}})$ this proves that then $D \in \text{safe}(R, \epsilon, L, \lambda, \omega)$.

3500 The proof is an application of Berge's maximum Theorem, Theorem D.3. Namely, define the function

$$3501 \quad f : \mathcal{R} \times \Pi \rightarrow \mathbb{R}, \quad f(R, \pi) := J_R(\pi) - \lambda \omega(\pi).$$

3502 Furthermore, define the correspondence $C : \mathcal{R} \rightrightarrows \Pi$ as the trivial map $C(R) = \Pi$. Let $f^* : \mathcal{R} \rightarrow \mathbb{R}$
3503 map a reward function to the value of a (λ, ω) -RLHF optimal policy, i.e., $f^*(R) := \max_{\pi \in \Pi} f(R, \pi)$.
3504 Define C^* as the corresponding argmax, i.e., $C^*(R) := \{\pi \mid f(R, \pi) = f^*(R)\}$. Assume on \mathcal{R}
3505 we have the standard Euclidean topology. Since ω is assumed continuous and by Proposition D.4
3506

3507 ⁴E.g., if $\pi_{\text{ref}}(a \mid s) > 0$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $\omega(\pi) := \mathbb{D}_{\text{KL}}(\pi \parallel \pi_{\text{ref}})$, then the minimum is given by
3508 π_{ref} .

also J is continuous, it follows that f is continuous. Thus, Theorem D.3 implies that C^* is upper hemicontinuous, see Definition D.2. The rest of the proof is simply an elaboration of why upper hemicontinuity of C^* gives the result.

Now, define the set

$$\mathcal{V} := \{\pi' \in \Pi \mid \text{Reg}^R(\pi') < \text{Reg}^R(\pi_{\text{ref}})\}.$$

Since the regret is a continuous function, this set is open. Now, let $\pi \in C^*(R)$ be (λ, ω) -RLHF optimal with respect to R . It follows

$$\begin{aligned} J_R(\pi) &= f(R, \pi) + \lambda\omega(\pi) \\ &> f(R, \pi_{\text{ref}}) + \lambda\omega(\pi_{\text{ref}}) \\ &= J_R(\pi_{\text{ref}}), \end{aligned}$$

where we used the optimality of π for f , that π_{ref} is not optimal for it, and that π_{ref} is the minimum of ω . So overall, this shows $C^*(R) \subseteq \mathcal{V}$.

Since C^* is upper hemicontinuous, this means there exists an open set $\mathcal{U} \subseteq \mathcal{R}$ with $R \in \mathcal{U}$ and such that for all $\hat{R} \in \mathcal{U}$, we have $C^*(\hat{R}) \subseteq \mathcal{V}$. Let $\epsilon > 0$ be so small that all reward functions \hat{R} with $\mathbb{E}_{(s,a) \sim D} \left[\frac{|\hat{R}(s,a) - R(s,a)|}{\text{range } R} \right] < \epsilon$ satisfy $\hat{R} \in \mathcal{U}$ — which exists since \mathcal{U} is open in the Euclidean topology. Then for all such \hat{R} and any policy $\hat{\pi}$ that is (λ, ω) -RLHF optimal wrt. \hat{R} , we by definition have

$$\hat{\pi} \in C^*(\hat{R}) \subseteq \mathcal{V},$$

and thus, by definition of \mathcal{V} , the desired regret property. This was to show. \square

Now, we show the same result, but with the choice distance instead of expected reward distance:

Theorem D.22. *Let $\lambda \in (0, \infty)$ be given and fixed. Assume we are given an MDP $\langle \mathcal{S}, \mathcal{A}, \tau, \mu_0, R, \gamma \rangle$, and a data distribution $D \in \mathcal{S} \times \mathcal{A}$ which assigns positive probability to all transitions, i.e., $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$, $D(s, a) > 0$. Let $\omega : \Pi \rightarrow \mathbb{R}$ be a continuous regularization function that has a reference policy π_{ref} as one of its minima. Assume that π_{ref} is not (λ, ω) -optimal for R and let $L = \text{Reg}^R(\pi_{\text{ref}})$. Then there exists $\epsilon > 0$ such that $D \in \text{safe}^{\text{DKL}}(R, \epsilon, L, \lambda, \omega)$.*

Proof. Let $\mathcal{G} := \mathbb{R}^{\Xi}$ be the vector space of return functions, which becomes a topological space when equipped with the infinity norm. Define the function

$$f : \mathcal{G} \times \Pi \rightarrow \mathbb{R}, \quad f(G, \pi) := J^G(\pi) - \lambda\omega(\pi),$$

where $J^G(\pi) := \mathbb{E}_{\xi \sim \pi} [G(\xi)]$ is the policy evaluation function of the return function G . f is continuous. Define the correspondence $C : \mathcal{G} \rightrightarrows \Pi$ as the trivial map $C(G) = \Pi$. Let $f^* : \mathcal{G} \rightarrow \mathbb{R}$ map a return function to the value of a (λ, ω) -optimal policy, i.e., $f^*(G) := \max_{\pi \in \Pi} f(G, \pi)$. Define C^* as the corresponding argmax. Then Theorem D.3 implies that C^* is upper hemicontinuous, see Definition D.2. As in the previous proof, the rest is an elaboration of why this gives the desired result.

Set G as the return function corresponding to R . Define

$$\mathcal{V} := \{\pi' \in \Pi \mid \text{Reg}^R(\pi') < L\}.$$

We now claim that $C^*(G) \subseteq \mathcal{V}$. Indeed, let $\pi \in C^*(G)$. Then

$$\begin{aligned} J^R(\pi) &= f(G, \pi) + \lambda\omega(\pi) \\ &> f(G, \pi_{\text{ref}}) + \lambda\omega(\pi_{\text{ref}}) \\ &= J^R(\pi_{\text{ref}}). \end{aligned}$$

Note that we used the optimality of π for f , that π_{ref} is not optimal for it, and also that π_{ref} minimizes ω by assumption. This shows $\text{Reg}^R(\pi) < \text{Reg}^R(\pi_{\text{ref}}) = L$, and thus the claim.

Since C^* is upper hemicontinuous and \mathcal{V} an open set, this implies that there exists $\sigma > 0$ such that for all $\hat{G} \in \mathcal{G}$ with $\|G - \hat{G}\|_{\infty} < \sigma$, we have $C^*(\hat{G}) \subseteq \mathcal{V}$.

Now, define $\epsilon := \epsilon(\sigma)$ as in Corollary D.20 and let \hat{R} be any reward function with $d_{\text{KL}}^D(R, \hat{R}) < \epsilon$. Then by that corollary, there exists $c \in \mathbb{R}$ such that $\|G - (\hat{G} - c)\|_{\infty} < \sigma$. Consequently, we have $C^*(\hat{G}) = C^*(\hat{G} - c) \subseteq \mathcal{V}$ by what we showed before, which shows the result. \square